

IOWA STATE UNIVERSITY

Digital Repository

Industrial and Manufacturing Systems Engineering
Publications

Industrial and Manufacturing Systems Engineering

2015

Multi-Agent Distributed Optimization via Inexact Consensus ADMM

Tsung-Hui Chang

National Taiwan University of Science and Technology

Mingyi Hong

Iowa State University, mingyi@iastate.edu

Xiangfeng Wang

East China Normal University

Follow this and additional works at: http://lib.dr.iastate.edu/imse_pubs



Part of the [Industrial Engineering Commons](#), and the [Systems Engineering Commons](#)

The complete bibliographic information for this item can be found at http://lib.dr.iastate.edu/imse_pubs/41. For information on how to cite this item, please visit <http://lib.dr.iastate.edu/howtocite.html>.

This Article is brought to you for free and open access by the Industrial and Manufacturing Systems Engineering at Iowa State University Digital Repository. It has been accepted for inclusion in Industrial and Manufacturing Systems Engineering Publications by an authorized administrator of Iowa State University Digital Repository. For more information, please contact digirep@iastate.edu.

Multi-Agent Distributed Optimization via Inexact Consensus ADMM

Abstract

Multi-agent distributed consensus optimization problems arise in many signal processing applications. Recently, the alternating direction method of multipliers (ADMM) has been used for solving this family of problems. ADMM based distributed optimization method is shown to have faster convergence rate compared with classic methods based on consensus subgradient, but can be computationally expensive, especially for problems with complicated structures or large dimensions. In this paper, we propose lowcomplexity algorithms that can reduce the overall computational cost of consensus ADMM by an order of magnitude for certain large-scale problems. Central to the proposed algorithms is the use of an inexact step for each ADMM update, which enables the agents to perform cheap computation at each iteration. Our convergence analyses show that the proposed methods converge well under some convexity assumptions. Numerical results show that the proposed algorithms offer considerably lower computational complexity than the standard ADMM based distributed optimization methods.

Keywords

distributed optimization, ADMM, consensus

Disciplines

Industrial Engineering | Systems Engineering

Comments

This is a manuscript of an article from IEEE Transactions on Signal Processing 63 (2015):482, doi: [10.1109/TSP.2014.2367458](https://doi.org/10.1109/TSP.2014.2367458). Posted with permission.

Multi-Agent Distributed Optimization via Inexact Consensus ADMM

Tsung-Hui Chang^{*}, Mingyi Hong[†] and Xiangfeng Wang[‡]

Abstract

Multi-agent distributed consensus optimization problems arise in many signal processing applications. Recently, the alternating direction method of multipliers (ADMM) has been used for solving this family of problems. ADMM based distributed optimization method is shown to have faster convergence rate compared with classic methods based on consensus subgradient, but can be computationally expensive, especially for problems with complicated structures or large dimensions. In this paper, we propose low-complexity algorithms that can reduce the overall computational cost of consensus ADMM by an order of magnitude for certain large-scale problems. Central to the proposed algorithms is the use of an inexact step for each ADMM update, which enables the agents to perform cheap computation at each iteration. Our convergence analyses show that the proposed methods converge well under some convexity assumptions. Numerical results show that the proposed algorithms offer considerably lower computational complexity than the standard ADMM based distributed optimization methods.

Keywords— Distributed optimization, ADMM, Consensus

EDICS: OPT-DOPT, MLR-DIST, NET-DISP, SPC-APPL.

The work of Tsung-Hui Chang is supported by National Science Council, Taiwan (R.O.C.), under Grant NSC 102-2221-E-011-005-MY3. Part of this work is presented in IEEE ICASSP 2014.

^{*}Tsung-Hui Chang is the corresponding author. Address: Department of Electronic and Computer Engineering, National Taiwan University of Science and Technology, Taipei 10607, Taiwan, (R.O.C.). E-mail: tsunghui.chang@ieee.org.

[†]Mingyi Hong is with Department of Electrical and Computer Engineering, University of Minnesota, Minneapolis, MN 55455, USA, E-mail: mhong@umn.edu

[‡]Xiangfeng Wang is with Shanghai Key Lab for Trustworthy Computing, Software Engineering Institute, East China Normal University, Shanghai, 200062, China, E-mail: xfwang@sei.ecnu.edu.cn

I. INTRODUCTION

We consider a network with multiple agents, for example a sensor network, a data cloud network or a communication network. The agents seek to collaborate to accomplish certain task. For example, distributed database servers may cooperate for data mining or for parameter learning in order to fully exploit the data collected from individual servers [1]. Another example arises from large-scale machine learning applications [2], where a computation task may be executed by collaborative microprocessors with individual memories and storage spaces [2]–[4]. Distributed optimization becomes favorable as it is not always efficient to pool all the local information for centralized computation, due to large size of problem dimension, a large amount of local data, energy constraints and/or privacy issues [5]–[8]. Many of the distributed optimization tasks, such as those described above, can be cast as an optimization problem of the following form

$$(P1) \quad \min_{\mathbf{y} \in \mathbb{R}^K} \sum_{i=1}^N \phi_i(\mathbf{y}) \quad (1)$$

where $\mathbf{y} \in \mathbb{R}^K$ is the decision variable and $\phi_i : \mathbb{R}^K \rightarrow \mathbb{R} \cup \{\infty\}$ is the cost function associated with agent i . Here the function ϕ_i is composed of a smooth component $f_i : \mathbb{R}^M \rightarrow \mathbb{R} \cup \{\infty\}$ (possibly with extended values) and a non-smooth component $g_i : \mathbb{R}^K \rightarrow \mathbb{R} \cup \{\infty\}$, i.e.,

$$\phi_i(\mathbf{y}) = f_i(\mathbf{A}_i \mathbf{y}) + g_i(\mathbf{y}), \quad (2)$$

where $\mathbf{A}_i \in \mathbb{R}^{M \times K}$ is some data matrix not necessarily of full rank. Such model is common in practice: the smooth component usually represents the cost function to be minimized, while the non-smooth component is often used as a regularization function [9] or an indicator function representing that \mathbf{y} is subject to a constraint set¹.

In the setting of distributed optimization, it is commonly assumed that each agent i only has knowledge about the local information f_i , g_i and \mathbf{A}_i . The challenge is to obtain, for each agent in the system, the optimal \mathbf{x} of (P1) using only local information and messages exchanged with neighbors [5]–[8].

¹For example, if $\mathbf{y} \in \mathcal{X} \subseteq \mathbb{R}^K$ for some set \mathcal{X} , then this can be implicitly included in the nonsmooth component g_i by letting [10, Section 5]

$$g_i(\mathbf{y}) = \begin{cases} 0 & \text{if } \mathbf{y} \in \mathcal{X} \\ \infty & \text{otherwise.} \end{cases} \quad (3)$$

In addition to (P1), another common problem formulation has the following form

$$(P2) \quad \min_{\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^K} \sum_{i=1}^N \phi_i(\mathbf{x}_i) \quad \text{s.t.} \quad \sum_{i=1}^N \mathbf{E}_i \mathbf{x}_i = \mathbf{q}, \quad (4)$$

where $\mathbf{E}_i \in \mathbb{R}^{M \times K}$, $\mathbf{q} \in \mathbb{R}^M$ and ϕ_i is given as in (2). Unlike (P1), in (P2), each agent i owns a local control variable² $\mathbf{x}_i \in \mathbb{R}^K$, and these variables are coupled together through the linear constraint. Examples of (P2) include the basis pursuit (BP) problem [11], [12], the network flow control problem [13] and interference management problem in communication networks [14]. To relate (P2) with (P1), let $\boldsymbol{\nu} \in \mathbb{R}^M$ be the Lagrange dual variable associated with the linear constraint $\sum_{i=1}^N \mathbf{E}_i \mathbf{x}_i = \mathbf{q}$. The Lagrange dual problem of (P2) can be equivalently written as

$$\min_{\boldsymbol{\nu} \in \mathbb{R}^M} \sum_{i=1}^N \left(\varphi_i(\boldsymbol{\nu}) + \frac{1}{N} \boldsymbol{\nu}^T \mathbf{q} \right) \quad (5)$$

where

$$\varphi_i(\boldsymbol{\nu}) = \max_{\mathbf{x}_i} \left\{ -\phi_i(\mathbf{x}_i) - \boldsymbol{\nu}^T \mathbf{E}_i \mathbf{x}_i \right\}, \quad i = 1, \dots, N. \quad (6)$$

Problem (5) thus has the same form as (P1). Given the optimal $\boldsymbol{\nu}$ of (5) and assuming that (P2) has a zero duality gap [15], each agent i can obtain the associated optimal variable \mathbf{x}_i by solving (6). Therefore, a distributed optimization method that can solve (P1) may also be used for (P2) through solving (5).

There is an extensive literature on distributed consensus optimization methods, such as the consensus subgradient methods; see [5], [6] and the recent developments in [7], [8], [16], [17]. The consensus subgradient methods are appealing owing to their simplicity and the ability to handle a wide range of problems. However, the convergence of the consensus subgradient methods are usually slow.

Recently, the alternating direction method of multipliers (ADMM) [10], [18] has become popular for solving problems with forms of (P1) and (P2) in a distributed fashion. In [14], distributed transmission designs for multi-cellular wireless communications were developed based on ADMM. In [19], several ADMM based distributed optimization algorithms were developed for solving the sparse LASSO problem [20]. In [12], using a different consensus formulation from [19] and assuming the availability of a certain coloring scheme for the graph, ADMM is applied to solving the BP problem [11] for both row partitioned and column partitioned data models [16]. In [21], the methodologies proposed in [12] are extended to handling a more general class of problems with forms of (P1) and (P2). In [22], a distributed ADMM with a sequential update rule is proposed; while in [23], the method is extended and can be implemented

²Here we let all \mathbf{x}_i 's have the same dimension without loss of generality.

asynchronously. The fast practical performance of ADMM is corroborated by its nice theoretical property. In particular, ADMM was found to converge linearly for a large class of problems [24], [25], meaning a certain optimality measure can decrease by a constant fraction in each iteration of the algorithm. In [26], [27], such fast convergence rate has also been built for distributed optimization.

It is important to note that existing ADMM based algorithms can be readily used to solve problems (P1) and (P2). For example, by applying the consensus formulation proposed in [19] and ADMM to (P1), a fully parallelized distributed optimization algorithm can be obtained (where the agents update their variables in a fully parallel manner), which we refer to as the consensus ADMM (C-ADMM). To solve (P2), the same consensus formulation and ADMM can be used on its Lagrange dual problem in (5), referred to as the dual consensus ADMM (DC-ADMM). The main drawback of these algorithms lies in the fact that each agent needs to repeatedly solve certain subproblems to *global optimality*. This can be computationally demanding, especially when the cost functions f_i 's have complicated structures or when the problem size is large [2]. If a low-accuracy suboptimal solution is used for these subproblems instead, the convergence is no longer guaranteed.

The main objective of this paper is to study algorithms that can significantly reduce the computational burden for the agents. In particular, we propose two algorithms, named the inexact consensus ADMM (IC-ADMM) and the inexact dual consensus ADMM (IDC-ADMM'), both of which allow the agents to perform a single proximal gradient (PG) step [28] at each iteration. The benefit of the proposed approach lies in the fact that the PG step is usually simple, especially when g_i 's are structured functions [9], [28]. Notably, the cheap iterations of the proposed algorithms is made possible by *inexactly* solving the subproblems arising in C-ADMM and DC-ADMM, in a way that is not known in the ADMM or consensus literature. For example, the proposed IC-ADMM approximates the smooth functions f_i 's in C-ADMM, which is very different from the known inexact ADMM methods [29], [30], where only the quadratic penalty is approximated (thus does not always result in cheap PG steps). We summarize our main contributions below.

- For (P1), we propose an IC-ADMM method for reducing the computational complexity of C-ADMM. Conditions for global convergence of IC-ADMM are analyzed. Moreover, we show that IC-ADMM converges linearly, under similar conditions as in [26].
- For (P2), we first propose a DC-ADMM method which can globally solve (P2) for any connected graph and convex ϕ_i 's. We further propose an IDC-ADMM method for reducing the computational burden of DC-ADMM. Conditions for global (linear) convergence are presented.

Numerical examples for solving distributed sparse logistic regression problems [31] will show that the proposed IC-ADMM and IDC-ADMM methods converge much faster than the consensus subgradient method [5]. Further, compared with the original C-ADMM and DC-ADMM, the proposed method can reduce the overall computational cost by an order of magnitude.

The paper is organized as follows. Section II presents the applications and assumptions. The C-ADMM and IC-ADMM are presented in Section III; while DC-ADMM and IDC-ADMM are presented in Section IV. Numerical results are given in Section V and conclusions are drawn in Section VI.

Notations: $\mathbf{A} \succeq \mathbf{0}$ ($\succ \mathbf{0}$) means that matrix \mathbf{A} is positive semidefinite (positive definite). \mathbf{I}_K is the $K \times K$ identity matrix; $\mathbf{1}_K$ is the K -dimensional all-one vector. $\|\mathbf{a}\|_2$ denotes the Euclidean norm of vector \mathbf{a} , and $\|\mathbf{z}\|_{\mathbf{A}}^2 \triangleq \mathbf{z}^T \mathbf{A} \mathbf{z}$ for some $\mathbf{A} \succeq \mathbf{0}$. Notation \otimes denotes the Kronecker product. $\text{diag}\{a_1, \dots, a_N\}$ is a diagonal matrix with the i th diagonal element being a_i ; while $\text{blkdiag}\{\mathbf{A}_1, \dots, \mathbf{A}_N\}$ is a block diagonal matrix with the i th diagonal block matrix being \mathbf{A}_i . $\lambda_{\max}(\mathbf{A})$ and $\lambda_{\min}(\mathbf{A})$ denote the maximum and minimum eigenvalues of matrix \mathbf{A} , respectively.

II. APPLICATIONS AND NETWORK MODEL

A. Application to Data Regression

As discussed in Section I, (P1) and (P2) arise in many problems in sensor networks, data networks and machine learning tasks. Here let us focus on the classical regression problems. We consider a general formulation that incorporates the LASSO [19] and logistic regression (LR) [31] as special instances. Let $\mathbf{A} = [\mathbf{A}_1^T, \dots, \mathbf{A}_N^T]^T \in \mathbb{R}^{NM \times K}$ denote a regression data matrix, where $\mathbf{A}_i \in \mathbb{R}^{M \times K}$ for all $i = 1, \dots, N$. For a row partitioned data (RPD) model [12, Fig. 1], [16], the distributed regression problem is given by

$$\min_{\mathbf{y} \in \mathbb{R}^K} \sum_{i=1}^N \Psi_i(\mathbf{y}; \mathbf{A}_i, \mathbf{b}_i), \quad (7)$$

where $\Psi_i(\mathbf{y}; \mathbf{A}_i, \mathbf{b}_i)$ is the cost function defined on the local regression data \mathbf{A}_i and a local response signal $\mathbf{b}_i \in \mathbb{R}^M$. For example, the LASSO problem has $\Psi_i(\mathbf{y}; \mathbf{A}_i, \mathbf{b}_i) = \|\mathbf{b}_i - \mathbf{A}_i \mathbf{y}\|_2^2 + g_i(\mathbf{y})$. Similarly, for the LR problem, one has

$$\Psi_i(\mathbf{y}; \mathbf{A}_i, \mathbf{b}_i) = \sum_{m=1}^M \log(1 + \exp(-b_{im} \mathbf{a}_{im}^T \mathbf{y})) + g_i(\mathbf{y}), \quad (8)$$

where $\mathbf{A}_i = [\mathbf{a}_{i1}, \dots, \mathbf{a}_{iM}]^T$ contains M training data vectors and $b_{im} \in \{\pm 1\}$ are binary labels for the training data. It is clear that (7) has the same form as (P1). Here, the non-smooth function g_i can be

1-norm for sparse regression, as well as mixture with an indicator functions specifying that \mathbf{y} is confined in certain constraint set.

On the other hand, let $\mathbf{E} = [\mathbf{E}_1, \dots, \mathbf{E}_N] \in \mathbb{R}^{M \times NK}$ denote a regression data matrix, where $\mathbf{E}_i \in \mathbb{R}^{M \times K}$ for all $i = 1, \dots, N$. Then, for the column partitioned data (CPD) model [12, Fig. 1], [16], the distributed regression problem is formulated as

$$\min_{\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^K} \sum_{i=1}^N \Psi_i(\mathbf{x}_i; \mathbf{E}_i, \mathbf{b}), \quad (9)$$

where the response signal \mathbf{b} is known to all agents while each agent i has a local regression variable $\mathbf{x}_i \in \mathbb{R}^K$ and local regression data matrix $\mathbf{E}_i = [\mathbf{e}_{i1}, \dots, \mathbf{e}_{iM}]^T \in \mathbb{R}^{M \times K}$. For example, the LR problem has

$$\Psi_i(\mathbf{x}_i; \mathbf{E}_i, \mathbf{b}) = \sum_{m=1}^M \log(1 + \exp(-b_m \sum_{i=1}^N \mathbf{e}_{im}^T \mathbf{x}_i)) + g_i(\mathbf{x}_i). \quad (10)$$

By introducing a slack variable $\mathbf{z} = [z_1, \dots, z_M]^T \triangleq \sum_{i=1}^N \mathbf{E}_i \mathbf{x}_i$, the CPD LR problem can be reformulated as

$$\begin{aligned} \min_{\substack{\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^K, \\ \mathbf{z} \in \mathbb{R}^M}} \left\{ \sum_{m=1}^M \log(1 + \exp(-b_m z_m)) + \sum_{i=1}^N g_i(\mathbf{x}_i) \right\} \\ \text{s.t. } \sum_{i=1}^N \mathbf{E}_i \mathbf{x}_i - \mathbf{z} = \mathbf{0}, \end{aligned} \quad (11)$$

which is an instance of (P2). In Section V, we will primarily test our algorithms on the RPD and CPD regression problems.

B. Network Model and Assumptions

Let an undirected graph \mathcal{G} denote a multi-agent network, which contains a node set $V = \{1, \dots, N\}$ and an edge set \mathcal{E} . An edge $(i, j) \in \mathcal{E}$ if and only if agent i and agent j can communicate with each other (i.e., neighbors). The edge set \mathcal{E} defines an adjacency matrix $\mathbf{W} \in \{0, 1\}^{N \times N}$, where $[\mathbf{W}]_{i,j} = 1$ if $(i, j) \in \mathcal{E}$ and $[\mathbf{W}]_{i,j} = 0$ otherwise. In addition, one can define an index subset $\mathcal{N}_i = \{j \in V \mid (i, j) \in \mathcal{E}\}$ for the neighbors of each agent i , and a degree matrix $\mathbf{D} = \text{diag}\{|\mathcal{N}_1|, \dots, |\mathcal{N}_N|\}$ (a diagonal matrix). With \mathbf{W} and \mathbf{D} , the Laplacian matrix of \mathcal{G} is given by $\mathbf{L} = \mathbf{D} - \mathbf{W}$ which is a positive semidefinite matrix (i.e., $\mathbf{L} \succeq \mathbf{0}$) and satisfies $\mathbf{L}\mathbf{1}_N = \mathbf{0}$ [32].

We make the following assumptions on \mathcal{G} and problems (P1) and (P2).

Assumption 1 *The undirected graph \mathcal{G} is connected.*

Assumption 1 implies that any two agents in the network can always influence each other in the long run. We also have the following assumptions on problems (P1) and (P2).

Assumption 2 (a) In (P1), the functions $\phi_i : \mathbb{R}^K \rightarrow \mathbb{R} \cup \{\infty\}$ are proper closed convex functions; at every \mathbf{y} for which both $f_i(\mathbf{A}_i \mathbf{y})$ and $g_i(\mathbf{y})$ are well defined and $\phi_i(\mathbf{y}) < \infty$, there exists at least one bounded subgradient $\partial\phi_i(\mathbf{y}) \in \mathbb{R}^K$ such that $\phi_i(\mathbf{x}) \geq \phi_i(\mathbf{y}) + (\partial\phi_i(\mathbf{y}))^T(\mathbf{x} - \mathbf{y}) \forall \mathbf{x} \in \mathbb{R}^K$. Moreover, the minimum of (P1) can be attained.

(b) In (P2), the functions $\phi_i : \mathbb{R}^K \rightarrow \mathbb{R} \cup \{\infty\}$ are proper closed convex functions; ϕ_i has at least one bounded subgradient at every \mathbf{x}_i for which both $f_i(\mathbf{A}_i \mathbf{x}_i)$ and $g_i(\mathbf{x}_i)$ are well defined and $\phi_i(\mathbf{x}_i) < \infty$; the minimum of (P2) is attained and so is its optimal dual value; moreover, strong duality holds for (P2).

Assumption 3 For all $i \in V$, the smooth function f_i in (2) is strongly convex, i.e., there exists some $\sigma_{f,i}^2 > 0$ such that

$$(\nabla f_i(\mathbf{y}) - \nabla f_i(\mathbf{x}))^T(\mathbf{y} - \mathbf{x}) \geq \sigma_{f,i}^2 \|\mathbf{y} - \mathbf{x}\|_2^2 \forall \mathbf{y}, \mathbf{x} \in \mathbb{R}^M.$$

Moreover, f_i has Lipschitz continuous gradients, i.e., there exists some $L_{f,i} > 0$ such that

$$\|\nabla f_i(\mathbf{y}) - \nabla f_i(\mathbf{x})\|_2 \leq L_{f,i} \|\mathbf{y} - \mathbf{x}\|_2 \quad \forall \mathbf{y}, \mathbf{x} \in \mathbb{R}^M. \quad (12)$$

Note that, even under Assumption 3, $\phi_i(\mathbf{x}) = f_i(\mathbf{A}_i \mathbf{x}) + g_i(\mathbf{x})$ is not necessarily strongly convex in \mathbf{x} since the matrix \mathbf{A}_i can be fat and rank deficient. Both the LASSO problem [19] and the LR function in (8) satisfy Assumption 3³.

III. DISTRIBUTED CONSENSUS ADMM

In Section III-A, we briefly review the original C-ADMM [19] for solving (P1). In Section III-B, we propose a computationally efficient inexact C-ADMM method.

A. Review of C-ADMM

Under Assumption 1, (P1) can be equivalently written as

$$\min_{\substack{\mathbf{y}_1, \dots, \mathbf{y}_N, \\ \{\mathbf{t}_{ij}\}}} \sum_{i=1}^N \phi_i(\mathbf{y}_i) \quad (13a)$$

$$\text{s.t. } \mathbf{y}_i = \mathbf{t}_{ij} \quad \forall j \in \mathcal{N}_i, i \in V, \quad (13b)$$

$$\mathbf{y}_j = \mathbf{t}_{ij} \quad \forall j \in \mathcal{N}_i, i \in V, \quad (13c)$$

³The logistic regression function $\log(1 + \exp(-x))$ is strongly convex given that x lies in a compact set.

where $\{\mathbf{t}_{ij}\}$ are slack variables. According to (13), each agent i can optimize its local function $f_i(\mathbf{A}_i \mathbf{y}_i) + g_i(\mathbf{y}_i)$ with respect to a local copy of \mathbf{y} , i.e., \mathbf{y}_i , under the consensus constraints in (13b) and (13c). In [19], ADMM is employed to solve (13) in a distributed manner. Let $\{\mathbf{u}_{ij}\}$ and $\{\mathbf{v}_{ij}\}$ denote the Lagrange dual variables associated with constraints (13b) and (13c), respectively. According to [19], ADMM leads to the following iterative updates at each iteration k :

$$\mathbf{u}_{ij}^{(k)} = \mathbf{u}_{ij}^{(k-1)} + \frac{c}{2}(\mathbf{y}_i^{(k-1)} - \mathbf{y}_j^{(k-1)}) \quad \forall j \in \mathcal{N}_i, i \in V, \quad (14a)$$

$$\mathbf{v}_{ij}^{(k)} = \mathbf{v}_{ij}^{(k-1)} + \frac{c}{2}(\mathbf{y}_j^{(k-1)} - \mathbf{y}_i^{(k-1)}) \quad \forall j \in \mathcal{N}_i, i \in V, \quad (14b)$$

$$\begin{aligned} \mathbf{y}_i^{(k)} = \arg \min_{\mathbf{y}_i} & \left\{ \phi_i(\mathbf{y}_i) + \sum_{j \in \mathcal{N}_i} (\mathbf{u}_{ij}^{(k)} + \mathbf{v}_{ji}^{(k)})^T \mathbf{y}_i \right. \\ & \left. + c \sum_{j \in \mathcal{N}_i} \left\| \mathbf{y}_i - \frac{\mathbf{y}_i^{(k-1)} + \mathbf{y}_j^{(k-1)}}{2} \right\|_2^2 \right\} \quad \forall i \in V, \end{aligned} \quad (14c)$$

where $c > 0$ is a penalty parameter and $\mathbf{u}_{ij}^{(0)} + \mathbf{v}_{ij}^{(0)} = \mathbf{0} \quad \forall i, j$. Note that variables $\{\mathbf{t}_{ij}^{(k)}\}$ are not shown in (14) as they can be expressed by variables $\{\mathbf{y}_i^{(k-1)}\}$; see [19] for the details.

The updates in (14) are useful for convergence analysis. For practical implementation, we define $\mathbf{p}_i^{(k)} \triangleq \sum_{j \in \mathcal{N}_i} (\mathbf{u}_{ij}^{(k)} + \mathbf{v}_{ji}^{(k)})$, $i \in V$. Then, (14) boils down to Algorithm 1.

Algorithm 1 C-ADMM for solving (P1)

1: **Given** initial variables $\mathbf{y}_i^{(0)} \in \mathbb{R}^K$ and $\mathbf{p}_i^{(0)} = \mathbf{0}$ for each agent i , $i \in V$. Set $k = 1$.

2: **repeat**

3: For all $i \in V$ (in parallel), $\mathbf{p}_i^{(k)} = \mathbf{p}_i^{(k-1)} + c \sum_{j \in \mathcal{N}_i} (\mathbf{y}_i^{(k-1)} - \mathbf{y}_j^{(k-1)})$,

$$\mathbf{y}_i^{(k)} = \arg \min_{\mathbf{y}_i} \left\{ f_i(\mathbf{A}_i \mathbf{y}_i) + g_i(\mathbf{y}_i) + \mathbf{y}_i^T \mathbf{p}_i^{(k)} + c \sum_{j \in \mathcal{N}_i} \left\| \mathbf{y}_i - \frac{\mathbf{y}_i^{(k-1)} + \mathbf{y}_j^{(k-1)}}{2} \right\|_2^2 \right\}. \quad (15)$$

4: **Set** $k = k + 1$.

5: **until** a predefined stopping criterion (e.g., a maximum iteration number) is satisfied.

It is important to note from Step 4 and Step 5 of Algorithm 1 that, except for the parameter c which has to be universally known, each agent i updates the variables $(\mathbf{y}_i^{(k)}, \mathbf{p}_i^{(k)})$ in a fully parallel manner, by only using the local function ϕ_i and messages $\{\mathbf{y}_j^{(k-1)}\}_{j \in \mathcal{N}_i}$, which come from its direct neighbors. It has been shown in [19] that, under Assumptions 1 and 2, C-ADMM is guaranteed to converge for any

$c > 0^4$:

$$\lim_{k \rightarrow \infty} \mathbf{y}_i^{(k)} = \mathbf{y}^*, \quad \lim_{k \rightarrow \infty} (\mathbf{u}_{ij}^{(k)}, \mathbf{v}_{ij}^{(k)}) = (\mathbf{u}_{ij}^*, \mathbf{v}_{ij}^*), \quad \forall j, i, \quad (16)$$

where $\mathbf{y}^* \triangleq \mathbf{y}_1^* = \dots = \mathbf{y}_N^*$ and $\{\mathbf{u}_{ij}^*, \mathbf{v}_{ij}^*\}$ denote a pair of optimal primal and dual solutions to problem (13), and \mathbf{y}^* is optimal to (P1). It is also shown that C-ADMM can converge linearly when ϕ_i 's are purely smooth (i.e., $g_i(\mathbf{y}_i) = 0 \ \forall i$) and strongly convex with respect to \mathbf{y}_i [26].

One key issue about C-ADMM is that the subproblem in (15) is not always easy to solve. For instance, for the LR function in (8), the associated subproblem (15) is given by

$$\begin{aligned} \mathbf{y}_i^{(k)} = \arg \min_{\mathbf{y}_i} & \left\{ \sum_{m=1}^M \log(1 + \exp(-b_{im} \mathbf{a}_{im}^T \mathbf{y}_i)) + g_i(\mathbf{y}_i) \right. \\ & \left. + \mathbf{y}_i^T \mathbf{p}_i^{(k)} + c \sum_{j \in \mathcal{N}_i} \left\| \mathbf{y}_i - \frac{\mathbf{y}_i^{(k-1)} + \mathbf{y}_j^{(k-1)}}{2} \right\|_2^2 \right\}. \end{aligned} \quad (17)$$

As seen, due to the complicated LR cost, problem (17) cannot yield simple solutions, and a numerical solver has to be employed. Clearly, obtaining a high-accuracy solution of (17) can be computationally expensive, especially when the problem dimension or the number of training data is large. While a low-accuracy solution to (17) can be adopted for complexity reduction, it may destroy the convergence behavior of C-ADMM, as will be shown in Section V.

B. Proposed Inexact C-ADMM

To reduce the complexity of C-ADMM, instead of solving subproblem (15) directly, we consider the following update:

$$\begin{aligned} \mathbf{y}_i^{(k)} = \arg \min_{\mathbf{y}_i} & \left\{ \nabla f_i(\mathbf{A}_i \mathbf{y}_i^{(k-1)})^T \mathbf{A}_i (\mathbf{y}_i - \mathbf{y}_i^{(k-1)}) \right. \\ & \left. + \frac{\beta_i}{2} \|\mathbf{y}_i - \mathbf{y}_i^{(k-1)}\|_2^2 + g_i(\mathbf{y}_i) + \mathbf{y}_i^T \mathbf{p}_i^{(k)} + c \sum_{j \in \mathcal{N}_i} \left\| \mathbf{y}_i - \frac{\mathbf{y}_i^{(k-1)} + \mathbf{y}_j^{(k-1)}}{2} \right\|_2^2 \right\}. \end{aligned} \quad (18)$$

In (18) we have replaced the smooth cost function $f_i(\mathbf{A}_i \mathbf{y}_i)$ in (15) with a proximal first-order approximation around $\mathbf{y}_i^{(k-1)}$:

$$\nabla f_i(\mathbf{A}_i \mathbf{y}_i^{(k-1)})^T \mathbf{A}_i (\mathbf{y}_i - \mathbf{y}_i^{(k-1)}) + \frac{\beta_i}{2} \|\mathbf{y}_i - \mathbf{y}_i^{(k-1)}\|_2^2,$$

where $\beta_i > 0$ is a penalty parameter of the proximal quadratic term. To obtain a concise representation of $\mathbf{y}_i^{(k)}$, let us define the *proximity operator* for the non-smooth function g_i at a given point $\mathbf{s} \in \mathbb{R}^K$ as

⁴In general, the parameter c is chosen empirically. Only for some special instance, optimal c may be analytically found; e.g., see [33].

[28]

$$\text{prox}_{g_i}^{\gamma_i}[\mathbf{s}] \triangleq \arg \min_{\mathbf{y}} \left\{ g_i(\mathbf{y}) + \frac{\gamma_i}{2} \|\mathbf{y} - \mathbf{s}\|_2^2 \right\}, \quad (19)$$

where $\gamma_i = \beta_i + 2c|\mathcal{N}_i|$. Clearly, using this definition, (18) can be expressed more compactly as

$$\begin{aligned} \mathbf{y}_i^{(k)} &= \arg \min_{\mathbf{y}_i} \left\{ g_i(\mathbf{y}_i) + \frac{\gamma_i}{2} \left\| \mathbf{y}_i - \frac{1}{\gamma_i} (\beta_i \mathbf{y}_i^{(k-1)} - \mathbf{p}_i^{(k)} \right. \right. \\ &\quad \left. \left. - \mathbf{A}_i^T \nabla f_i(\mathbf{A}_i \mathbf{y}_i^{(k-1)}) + c \sum_{j \in \mathcal{N}_i} (\mathbf{y}_i^{(k-1)} + \mathbf{y}_j^{(k-1)}) \right\|_2^2 \right\} \\ &= \text{prox}_{g_i}^{\gamma_i} \left[\frac{1}{\gamma_i} \left(\beta_i \mathbf{y}_i^{(k-1)} - \mathbf{p}_i^{(k)} - \mathbf{A}_i^T \nabla f_i(\mathbf{A}_i \mathbf{y}_i^{(k-1)}) \right. \right. \\ &\quad \left. \left. + c \sum_{j \in \mathcal{N}_i} (\mathbf{y}_i^{(k-1)} + \mathbf{y}_j^{(k-1)}) \right) \right], \end{aligned} \quad (20)$$

which is a proximal gradient (PG) update.

The PG updates like (20) often admit closed-form expression, especially when g_i 's are functions including the ℓ_1 norm, Euclidean norm, infinity norm and matrix nuclear norm [34]. For example, when $g_i(\mathbf{y}) = \|\mathbf{y}\|_1$, (19) has a closed-form solution known as the soft thresholding operator [28], [34]:

$$\mathcal{S} \left[\mathbf{s}, \frac{1}{\gamma_i} \right] = \left(\mathbf{s} - \frac{1}{\gamma_i} \mathbf{1}_K \right)^+ + \left(-\mathbf{s} - \frac{1}{\gamma_i} \mathbf{1}_K \right)^+, \quad (21)$$

where $(x)^+ \triangleq \max\{x, 0\}$. The IC-ADMM is presented in Algorithm 2.

Algorithm 2 Proposed IC-ADMM for solving (P1)

1: **Given** initial variables $\mathbf{y}_i^{(0)} \in \mathbb{R}^K$ and $\mathbf{p}_i^{(0)} = \mathbf{0}$ for each agent i , $i \in V$. Set $k = 1$.

2: **repeat**

3: For all $i \in V$ (in parallel),

$$\begin{aligned} \mathbf{p}_i^{(k)} &= \mathbf{p}_i^{(k-1)} + c \sum_{j \in \mathcal{N}_i} (\mathbf{y}_i^{(k-1)} - \mathbf{y}_j^{(k-1)}), \\ \mathbf{y}_i^{(k)} &= \text{prox}_{g_i}^{\gamma_i} \left[\frac{1}{\gamma_i} \left(\beta_i \mathbf{y}_i^{(k-1)} - \mathbf{A}_i^T \nabla f_i(\mathbf{A}_i \mathbf{y}_i^{(k-1)}) \right. \right. \\ &\quad \left. \left. - \mathbf{p}_i^{(k)} + c \sum_{j \in \mathcal{N}_i} (\mathbf{y}_i^{(k-1)} + \mathbf{y}_j^{(k-1)}) \right) \right]. \end{aligned} \quad (22)$$

4: **Set** $k = k + 1$.

5: **until** a predefined stopping criterion (e.g., a maximum iteration number) is satisfied.

Although the idea of “inexact ADMM” is not new, our approach is significantly different from the existing methods [29], [30], where the inexact update is obtained by approximating the quadratic penalization term only. It can be seen that problem (17) is still difficult to solve even the inexact update

in [29], [30] is applied. Two notable exceptions are the algorithms proposed in [35] and [36] where the cost function is also linearized. However, an additional back substitution step and two extragradient steps are required in [35] and [36], respectively, which is not suited for distributed optimization.

The convergence properties of IC-ADMM is characterized by the following theorem.

Theorem 1 *Suppose that Assumptions 1, 2(a) and 3 hold. Let*

$$\beta_i > \frac{L_{f,i}^2}{\sigma_{f,i}^2} \lambda_{\max}(\mathbf{A}_i^T \mathbf{A}_i) - c \lambda_{\min}(\mathbf{D} + \mathbf{W}) > 0 \quad \forall i \in V, \quad (23)$$

and let $\mathbf{y}^ \triangleq \mathbf{y}_1^* = \dots = \mathbf{y}_N^*$ and $\{\mathbf{u}_{ij}^*, \mathbf{v}_{ij}^*\}$ denote a pair of optimal primal and dual solutions to problem (13) (i.e., (P1)).*

(a) For Algorithm 2, $\mathbf{y}_1^{(k)}, \dots, \mathbf{y}_N^{(k)}$ converge to a common point \mathbf{y}^ .*

(b) If $\phi_i(\mathbf{y}) = f_i(\mathbf{A}_i \mathbf{y})$, where \mathbf{A}_i has full column rank, for all $i \in V$, then we have

$$\lim_{k \rightarrow \infty} \|\mathbf{y}^{(k)} - \mathbf{1}_N \otimes \mathbf{y}^*\|_{\frac{1}{2}\mathbf{G} + \alpha\mathbf{M}}^2 + \frac{1}{c} \|\mathbf{u}^{(k+1)} - \mathbf{u}^*\|_2^2 = 0 \text{ linearly,}$$

where $\mathbf{y}^{(k)} = [(\mathbf{y}_1^{(k)})^T, \dots, (\mathbf{y}_N^{(k)})^T]^T$; $\mathbf{u}_i^{(k)} \in \mathbb{R}^{K|\mathcal{N}_i|}$ (\mathbf{u}_i^) is a vector that stacks $\mathbf{u}_{ij}^{(k)}$ (\mathbf{u}_{ij}^*) $\forall j \in \mathcal{N}_i$; $\mathbf{u}^{(k)} \in \mathbb{R}^{K \sum_{i=1}^N |\mathcal{N}_i|}$ (\mathbf{u}^*) stacks $\mathbf{u}_i^{(k)}$ (\mathbf{u}_i^*) $\forall i = 1, \dots, N$. and*

$$\mathbf{G} \triangleq \mathbf{D}_\beta + c((\mathbf{D} + \mathbf{W}) \otimes \mathbf{I}_K) \succ \mathbf{0}, \quad (24)$$

$$\mathbf{M} \triangleq \tilde{\mathbf{A}}^T (\mathbf{D}_{\sigma_f} - \frac{1}{2} \mathbf{D}_\rho) \tilde{\mathbf{A}} \succ \mathbf{0}, \quad (25)$$

for some $0 < \alpha < 1$ and $\rho > 0$. Here, $\tilde{\mathbf{A}} = \text{blkdiag}\{\mathbf{A}_1, \dots, \mathbf{A}_N\}$; $\mathbf{D}_\beta = \text{diag}\{\beta_1, \dots, \beta_N\} \otimes \mathbf{I}_K$; $\mathbf{D}_{\sigma_f} = \text{diag}\{\sigma_{f,1}^2, \dots, \sigma_{f,N}^2\} \otimes \mathbf{I}_K$; and $\mathbf{D}_\rho = \text{diag}\{\rho_1, \dots, \rho_N\} \otimes \mathbf{I}_K$.

The proof is presented in Appendix A. Theorem 1 implies that, given sufficiently large β_i 's, IC-ADMM not only achieves consensus and optimality, but also converges linearly provided that ϕ_i is purely smooth and strongly convex. Note that, to ensure (23), the global knowledge of $\lambda_{\min}(\mathbf{D} + \mathbf{W})$ is required by all agents. As a parallel work, we should mention that a concurrent result similar as Theorem 1(b) is presented in [37].

Remark 1 We remark that the convergence condition in (23) depends on the network topology. Let $\mathbf{L} = \mathbf{D} - \mathbf{W}$ denote the Laplacian matrix of \mathcal{G} . Then $\mathbf{D} + \mathbf{W} = 2\mathbf{D} - \mathbf{L}$. By the graph theory [32], the normalized Laplacian matrix, i.e., $\tilde{\mathbf{L}} = \mathbf{D}^{-\frac{1}{2}} \mathbf{L} \mathbf{D}^{-\frac{1}{2}}$, must have $\lambda_{\max}(\tilde{\mathbf{L}}) \leq 2$. Further, $\lambda_{\max}(\tilde{\mathbf{L}}) < 2$ if and only if the connected graph \mathcal{G} is not bipartite. Thus, we have $\lambda_{\min}(\mathbf{D} + \mathbf{W}) = \lambda_{\min}(\mathbf{D}^{\frac{1}{2}}(2\mathbf{I}_N - \tilde{\mathbf{L}})\mathbf{D}^{\frac{1}{2}}) \geq 0$, and $\lambda_{\min}(\mathbf{D} + \mathbf{W}) > 0$ whenever \mathcal{G} is non-bipartite.

IV. DISTRIBUTED DUAL CONSENSUS ADMM

In this section, we turn the focus to (P2). In Section IV-A, we present a DC-ADMM method for solving (P2). In Section IV-B, an inexact DC-ADMM method is proposed.

A. Proposed DC-ADMM

The DC-ADMM is obtained by applying the C-ADMM (Algorithm 1) to problem (5) which is equivalent to the Lagrange dual of (P2). Firstly, similar to (13), we write problem (5) as

$$\min_{\substack{\boldsymbol{\nu}_1, \dots, \boldsymbol{\nu}_N \\ \{\mathbf{t}_{ij}\}}} \sum_{i=1}^N \left(\varphi_i(\boldsymbol{\nu}_i) + \frac{1}{N} \boldsymbol{\nu}_i^T \mathbf{q} \right) \quad (26a)$$

$$\text{s.t. } \boldsymbol{\nu}_i = \mathbf{t}_{ij}, \boldsymbol{\nu}_j = \mathbf{t}_{ij} \quad \forall j \in \mathcal{N}_i, i \in V, \quad (26b)$$

where $\boldsymbol{\nu}_i \in \mathbb{R}^M$ is the i th agent's local copy of the dual variable $\boldsymbol{\nu}$ and φ_i is given in (6). Following a similar argument as in deriving Algorithm 1, we obtain the following update steps at each iteration k

$$\mathbf{p}_i^{(k)} = \mathbf{p}_i^{(k-1)} + c \sum_{j \in \mathcal{N}_i} (\boldsymbol{\nu}_i^{(k-1)} - \boldsymbol{\nu}_j^{(k-1)}), \quad (27a)$$

$$\begin{aligned} \boldsymbol{\nu}_i^{(k)} = \arg \min_{\boldsymbol{\nu}_i \in \mathbb{R}^M} & \left\{ \varphi_i(\boldsymbol{\nu}_i) + \frac{1}{N} \boldsymbol{\nu}_i^T \mathbf{q} + \boldsymbol{\nu}_i^T \mathbf{p}_i^{(k)} \right. \\ & \left. + c \sum_{j \in \mathcal{N}_i} \left\| \boldsymbol{\nu}_i - \frac{\boldsymbol{\nu}_i^{(k-1)} + \boldsymbol{\nu}_j^{(k-1)}}{2} \right\|_2^2 \right\} \quad \forall i \in V, \end{aligned} \quad (27b)$$

where, with a slight abuse of notation,

$$\mathbf{p}_i^{(k)} = \sum_{j \in \mathcal{N}_i} (\mathbf{u}_{ij}^{(k)} + \mathbf{v}_{ji}^{(k)}), \quad (28)$$

in which $\{\mathbf{u}_{ij}\}$ and $\{\mathbf{v}_{ij}\}$ are dual variables associated with the two constraints in (26b) and are updated in a similar fashion as in (14a) and (14b), i.e.,

$$\mathbf{u}_{ij}^{(k)} = \mathbf{u}_{ij}^{(k-1)} + \frac{c}{2} (\boldsymbol{\nu}_i^{(k-1)} - \boldsymbol{\nu}_j^{(k-1)}) \quad \forall j \in \mathcal{N}_i, i \in V, \quad (29a)$$

$$\mathbf{v}_{ij}^{(k)} = \mathbf{v}_{ij}^{(k-1)} + \frac{c}{2} (\boldsymbol{\nu}_j^{(k-1)} - \boldsymbol{\nu}_i^{(k-1)}) \quad \forall j \in \mathcal{N}_i, i \in V. \quad (29b)$$

In general, subproblem (27b) is not easy to handle because φ_i is implicit and (27b) is in fact a min-max optimization problem given by

$$\begin{aligned} \boldsymbol{\nu}_i^{(k)} = \arg \min_{\boldsymbol{\nu}_i} \max_{\mathbf{x}_i} & \left\{ -\phi_i(\mathbf{x}_i) - \boldsymbol{\nu}_i^T \mathbf{E}_i \mathbf{x}_i + \frac{1}{N} \boldsymbol{\nu}_i^T \mathbf{q} \right. \\ & \left. + \boldsymbol{\nu}_i^T \mathbf{p}_i^{(k)} + c \sum_{j \in \mathcal{N}_i} \left\| \boldsymbol{\nu}_i - \frac{\boldsymbol{\nu}_i^{(k-1)} + \boldsymbol{\nu}_j^{(k-1)}}{2} \right\|_2^2 \right\}. \end{aligned} \quad (30)$$

Fortunately, since the objective function in (30) is convex in $\boldsymbol{\nu}_i$ for any \mathbf{x}_i and is concave in \mathbf{x}_i for any $\boldsymbol{\nu}_i$, the minimax theorem [38, Proposition 2.6.2] can be applied so that the min-max problem (30) can

be equivalently solved by considering its max-min counterpart and saddle point exists. Specifically, the max-min counterpart of (30) is given by

$$\max_{\mathbf{x}_i} \min_{\boldsymbol{\nu}_i} \left\{ -\phi_i(\mathbf{x}_i) - \boldsymbol{\nu}_i^T \mathbf{E}_i \mathbf{x}_i + \frac{1}{N} \boldsymbol{\nu}_i^T \mathbf{q} + \boldsymbol{\nu}_i^T \mathbf{p}_i^{(k)} + c \sum_{j \in \mathcal{N}_i} \left\| \boldsymbol{\nu}_i - \frac{\boldsymbol{\nu}_i^{(k-1)} + \boldsymbol{\nu}_j^{(k-1)}}{2} \right\|_2^2 \right\} \quad (31)$$

$$\begin{aligned} = & \max_{\mathbf{x}_i} \min_{\boldsymbol{\nu}_i} \left\{ -\phi_i(\mathbf{x}_i) + (c|\mathcal{N}_i|) \left\| \boldsymbol{\nu}_i - \frac{1}{2|\mathcal{N}_i|} \left[\sum_{j \in \mathcal{N}_i} (\boldsymbol{\nu}_i^{(k-1)} + \boldsymbol{\nu}_j^{(k-1)}) - \frac{1}{c} \mathbf{p}_i^{(k)} + \frac{1}{c} (\mathbf{E}_i \mathbf{x}_i - \frac{1}{N} \mathbf{q}) \right] \right\|_2^2 \right. \\ & \left. - \frac{c}{4|\mathcal{N}_i|} \left\| \frac{1}{c} (\mathbf{E}_i \mathbf{x}_i - \frac{1}{N} \mathbf{q}) - \frac{1}{c} \mathbf{p}_i^{(k)} + \sum_{j \in \mathcal{N}_i} (\boldsymbol{\nu}_i^{(k-1)} + \boldsymbol{\nu}_j^{(k-1)}) \right\|_2^2 \right\} \end{aligned} \quad (32)$$

where the equality is obtained by completing the quadratic term of $\boldsymbol{\nu}_i$. Let $\mathbf{x}_i^{(k)}$ be an inner maximizer of (30) so that $(\boldsymbol{\nu}_i^{(k)}, \mathbf{x}_i^{(k)})$ is a saddle point of (30). Then, $(\mathbf{x}_i^{(k)}, \boldsymbol{\nu}_i^{(k)})$ is a pair of outer-inner solution to (31) and (32) [38, Proposition 2.6.1]. From (32), the inner minimizer $\boldsymbol{\nu}_i^{(k)}$ can be uniquely determined by

$$\begin{aligned} \boldsymbol{\nu}_i^{(k)} = & \frac{1}{2|\mathcal{N}_i|} \left[\sum_{j \in \mathcal{N}_i} (\boldsymbol{\nu}_i^{(k-1)} + \boldsymbol{\nu}_j^{(k-1)}) - \frac{1}{c} \mathbf{p}_i^{(k)} \right. \\ & \left. + \frac{1}{c} (\mathbf{E}_i \mathbf{x}_i^{(k)} - \frac{1}{N} \mathbf{q}) \right], \end{aligned} \quad (33)$$

and that the outer maximizer is given by

$$\begin{aligned} \mathbf{x}_i^{(k)} = & \arg \min_{\mathbf{x}_i} \left\{ \phi_i(\mathbf{x}_i) + \frac{c}{4|\mathcal{N}_i|} \left\| \frac{1}{c} (\mathbf{E}_i \mathbf{x}_i - \frac{1}{N} \mathbf{q}) - \frac{1}{c} \mathbf{p}_i^{(k)} + \sum_{j \in \mathcal{N}_i} (\boldsymbol{\nu}_i^{(k-1)} + \boldsymbol{\nu}_j^{(k-1)}) \right\|_2^2 \right\}. \end{aligned} \quad (34)$$

As a result, the min-max subproblem (27b) can actually be obtained by first solving the subproblem (34) with respect to the primal variable \mathbf{x}_i followed by evaluating $\boldsymbol{\nu}_i^{(k)}$ using the close-form in (33). The proposed DC-ADMM is summarized in Algorithm 3.

Interestingly, while DC-ADMM handles the equivalent dual problem in (5), it directly yields primal optimal solution of (P2), as we state in the following theorem.

Theorem 2 Suppose that Assumptions 1 and 2(b) hold. Then $(\boldsymbol{\nu}_1^{(k)}, \dots, \boldsymbol{\nu}_N^{(k)})$ converges to a common point $\boldsymbol{\nu}^*$, which is optimal to the dual problem (5). Moreover, any limit point of $(\mathbf{x}_1^{(k)}, \dots, \mathbf{x}_N^{(k)})$ is primal optimal to (P2).

Algorithm 3 Proposed DC-ADMM for solving (P2)

1: **Given** initial variables $\mathbf{x}_i^{(0)} \in \mathbb{R}^K$, $\boldsymbol{\nu}_i^{(0)} \in \mathbb{R}^M$ and $\mathbf{p}_i^{(0)} = \mathbf{0}$ for each agent i , $i \in V$. Set $k = 1$.

2: **repeat**

3: For all $i \in V$ (in parallel),

$$\begin{aligned} \mathbf{p}_i^{(k)} &= \mathbf{p}_i^{(k-1)} + c \sum_{j \in \mathcal{N}_i} (\boldsymbol{\nu}_i^{(k-1)} - \boldsymbol{\nu}_j^{(k-1)}), \\ \mathbf{x}_i^{(k)} &= \arg \min_{\mathbf{x}_i} \left\{ \phi_i(\mathbf{x}_i) + \frac{c}{4|\mathcal{N}_i|} \left\| \frac{1}{c} (\mathbf{E}_i \mathbf{x}_i - \frac{1}{N} \mathbf{q}) \right. \right. \\ &\quad \left. \left. - \frac{1}{c} \mathbf{p}_i^{(k)} + \sum_{j \in \mathcal{N}_i} (\boldsymbol{\nu}_i^{(k-1)} + \boldsymbol{\nu}_j^{(k-1)}) \right\|_2^2 \right\}, \end{aligned} \quad (35)$$

$$\begin{aligned} \boldsymbol{\nu}_i^{(k)} &= \frac{1}{2|\mathcal{N}_i|} \left(\sum_{j \in \mathcal{N}_i} (\boldsymbol{\nu}_i^{(k-1)} + \boldsymbol{\nu}_j^{(k-1)}) - \frac{1}{c} \mathbf{p}_i^{(k)} \right. \\ &\quad \left. + \frac{1}{c} (\mathbf{E}_i \mathbf{x}_i^{(k)} - \frac{1}{N} \mathbf{q}) \right). \end{aligned} \quad (36)$$

4: **Set** $k = k + 1$.

5: **until** a predefined stopping criterion is satisfied.

Proof: Since DC-ADMM is a direct application of C-ADMM to the dual problem (5), it follows from [19] that as $k \rightarrow \infty$,

$$\boldsymbol{\nu}_i^{(k)} \rightarrow \boldsymbol{\nu}^*, \quad \boldsymbol{\nu}_i^{(k)} - \boldsymbol{\nu}_j^{(k)} \rightarrow \mathbf{0} \quad \forall j \in \mathcal{N}_i, \quad i \in V. \quad (37)$$

What remains is to show that any limit point of $(\mathbf{x}_1^{(k)}, \dots, \mathbf{x}_N^{(k)})$ is asymptotically optimal to (P2), i.e., as $k \rightarrow \infty$,

$$\partial \phi_i(\mathbf{x}_i^{(k)}) + \mathbf{E}_i^T \boldsymbol{\nu}^* \rightarrow \mathbf{0} \quad \forall i \in V, \quad (38)$$

$$\sum_{i=1}^N \mathbf{E}_i \mathbf{x}_i^{(k)} - \mathbf{q} \rightarrow \mathbf{0}. \quad (39)$$

To show (38), consider the optimality condition of (34), i.e.,

$$\begin{aligned} \mathbf{0} &= \partial \phi_i(\mathbf{x}_i^{(k)}) + \frac{1}{2|\mathcal{N}_i|} \mathbf{E}_i^T \left(\frac{1}{c} (\mathbf{E}_i \mathbf{x}_i^{(k)} - \frac{1}{N} \mathbf{q}) \right. \\ &\quad \left. - \frac{1}{c} \mathbf{p}_i^{(k)} + \sum_{j \in \mathcal{N}_i} (\boldsymbol{\nu}_i^{(k-1)} + \boldsymbol{\nu}_j^{(k-1)}) \right) \\ &= \partial \phi_i(\mathbf{x}_i^{(k)}) + \mathbf{E}_i^T \boldsymbol{\nu}_i^{(k)}, \end{aligned} \quad (40)$$

where the second equality is obtained by (33). Since (40) holds for all k and $\boldsymbol{\nu}_i^{(k)} \rightarrow \boldsymbol{\nu}^*$ by (37), (38) is true when $k \rightarrow \infty$.

To show (39), rewrite (33) as follows

$$\begin{aligned}
\mathbf{0} &= -(\mathbf{E}_i \mathbf{x}_i^{(k)} - \frac{1}{N} \mathbf{q}) + 2c \sum_{j \in \mathcal{N}_i} (\boldsymbol{\nu}_i^{(k)} - \frac{\boldsymbol{\nu}_i^{(k)} + \boldsymbol{\nu}_j^{(k)}}{2}) \\
&\quad + \mathbf{p}_i^{(k)} + c \sum_{j \in \mathcal{N}_i} (\boldsymbol{\nu}_i^{(k)} + \boldsymbol{\nu}_j^{(k)} - \boldsymbol{\nu}_i^{(k-1)} - \boldsymbol{\nu}_j^{(k-1)}) \\
&= -(\mathbf{E}_i \mathbf{x}_i^{(k)} - \frac{1}{N} \mathbf{q}) + \mathbf{p}_i^{(k+1)} \\
&\quad + c \sum_{j \in \mathcal{N}_i} (\boldsymbol{\nu}_i^{(k)} + \boldsymbol{\nu}_j^{(k)} - \boldsymbol{\nu}_i^{(k-1)} - \boldsymbol{\nu}_j^{(k-1)}), \tag{41}
\end{aligned}$$

where the last equality is obtained by (28) and (29). Upon summing (41) for $i = 1, \dots, N$, and by the fact that

$$\sum_{i=1}^N \mathbf{p}_i^{(k)} = \sum_{i=1}^N \sum_{j \in \mathcal{N}_i} (\mathbf{u}_{ij}^{(k)} + \mathbf{v}_{ji}^{(k)}) = \mathbf{0}$$

(by applying (A.13) and (A.14) in Appendix A), we can obtain

$$\sum_{i=1}^N \mathbf{E}_i \mathbf{x}_i^{(k)} - \mathbf{q} = c \sum_{i=1}^N \sum_{j \in \mathcal{N}_i} (\boldsymbol{\nu}_i^{(k)} + \boldsymbol{\nu}_j^{(k)} - \boldsymbol{\nu}_i^{(k-1)} - \boldsymbol{\nu}_j^{(k-1)}). \tag{42}$$

Note that $\boldsymbol{\nu}_i^{(k)} - \boldsymbol{\nu}_i^{(k-1)} \rightarrow \mathbf{0} \ \forall i \in V$ as inferred from $\boldsymbol{\nu}_i^{(k)} \rightarrow \boldsymbol{\nu}^* \ \forall i \in V$ in (37). By applying this fact to (42), we obtain that (39) is true as $k \rightarrow \infty$. \blacksquare

Interestingly, from (42), one observes that the primal feasibility of $(\mathbf{x}_1^{(k)}, \dots, \mathbf{x}_N^{(k)})$ to (P2) depends on the agents' consensus on the dual variable $\boldsymbol{\nu}$.

We remark that Algorithm 3 is different from the D-ADMM algorithm in [12, Algorithm 3]. Firstly, Algorithm 3 can be implemented in a fully parallel manner; secondly, Algorithm 3 does not involve solving a min-max subproblem at each iteration; thirdly, convergence of Algorithm 3 can be achieved without the assumption that the graph \mathcal{G} is bipartite.

B. Proposed Inexact DC-ADMM

In this subsection, we propose an inexact version of DC-ADMM, referred to as the IDC-ADMM. In view of the fact that solving the subproblem in (35) can be expensive, we consider an inexact update of $\mathbf{x}_i^{(k)}$. Specifically, since a non-trivial \mathbf{E}_i can also complicate the solution⁵, we propose to approximate both $f_i(\mathbf{A}_i \mathbf{x}_i)$ and the quadratic term $\frac{c}{4|\mathcal{N}_i|} \|\frac{1}{c}(\mathbf{E}_i \mathbf{x}_i - \frac{1}{N} \mathbf{q}) - \frac{1}{c} \mathbf{p}_i^{(k)} + \sum_{j \in \mathcal{N}_i} (\boldsymbol{\nu}_i^{(k-1)} + \boldsymbol{\nu}_j^{(k-1)})\|_2^2$ in (35)

⁵When \mathbf{E}_i has orthogonal columns (e.g., $\mathbf{E}_i^T \mathbf{E}_i = \alpha \mathbf{I}_K$ for some $\alpha \in \mathbb{R}$), then it may not be necessary to approximate the quadratic term.

by a proximal first-order approximation around $\mathbf{x}_i^{(k-1)}$; this leads to the following update

$$\begin{aligned} \mathbf{x}_i^{(k)} = \arg \min_{\mathbf{x}_i} & \left\{ \left[\mathbf{A}_i^T \nabla f_i(\mathbf{A}_i \mathbf{x}_i^{(k-1)}) + \frac{1}{2|\mathcal{N}_i|} \mathbf{E}_i^T \left(\frac{1}{c} (\mathbf{E}_i \mathbf{x}_i^{(k-1)} \right. \right. \right. \\ & \left. \left. \left. - \frac{1}{N} \mathbf{q} \right) - \frac{1}{c} \mathbf{p}_i^{(k)} + \sum_{j \in \mathcal{N}_i} (\boldsymbol{\nu}_i^{(k-1)} + \boldsymbol{\nu}_j^{(k-1)}) \right) \right]^T (\mathbf{x}_i - \mathbf{x}_i^{(k-1)}) \\ & \left. + \frac{\beta_i}{2} \|\mathbf{x}_i - \mathbf{x}_i^{(k-1)}\|_2^2 + g_i(\mathbf{x}_i) \right\}, \end{aligned} \quad (43)$$

where, with a slight abuse of notation, $\beta_i > 0$ is a penalty parameter. By (19), equation (43) can be further written as the following PG update

$$\begin{aligned} \mathbf{x}_i^{(k)} = \arg \min_{\mathbf{x}_i} & \left\{ \frac{\beta_i}{2} \left\| \mathbf{x}_i - \left[\mathbf{x}_i^{(k-1)} - \frac{1}{\beta_i} \mathbf{A}_i^T \nabla f_i(\mathbf{A}_i \mathbf{x}_i^{(k-1)}) \right. \right. \right. \\ & \left. \left. \left. - \frac{1}{2\beta_i |\mathcal{N}_i|} \mathbf{E}_i^T \left(\frac{1}{c} (\mathbf{E}_i \mathbf{x}_i^{(k-1)} - \frac{1}{N} \mathbf{q}) - \frac{1}{c} \mathbf{p}_i^{(k)} \right. \right. \right. \\ & \left. \left. \left. + \sum_{j \in \mathcal{N}_i} (\boldsymbol{\nu}_i^{(k-1)} + \boldsymbol{\nu}_j^{(k-1)}) \right) \right] \right\|_2^2 + g_i(\mathbf{x}_i) \right\} \\ = \text{prox}_{g_i}^{\beta_i} & \left[\mathbf{x}_i^{(k-1)} - \frac{1}{\beta_i} \mathbf{A}_i^T \nabla f_i(\mathbf{A}_i \mathbf{x}_i^{(k-1)}) \right. \\ & \left. - \frac{1}{2\beta_i |\mathcal{N}_i|} \mathbf{E}_i^T \left(\frac{1}{c} (\mathbf{E}_i \mathbf{x}_i^{(k-1)} - \frac{1}{N} \mathbf{q}) - \frac{1}{c} \mathbf{p}_i^{(k)} \right. \right. \\ & \left. \left. + \sum_{j \in \mathcal{N}_i} (\boldsymbol{\nu}_i^{(k-1)} + \boldsymbol{\nu}_j^{(k-1)}) \right) \right]. \end{aligned} \quad (44)$$

We summarize the proposed IDC-ADMM in Algorithm 4.

The convergence property of IDC-ADMM is stated below.

Theorem 3 Suppose that Assumptions 1, 2(b) and 3 hold and

$$\beta_i > \lambda_{\max} \left(\frac{L_{f,i}^2}{\sigma_{f,i}^2} \mathbf{A}_i^T \mathbf{A}_i + \frac{1}{2|\mathcal{N}_i|c} \mathbf{E}_i^T \mathbf{E}_i \right) \quad \forall i \in V. \quad (47)$$

Let $\mathbf{x}^* = [(\mathbf{x}_1^*)^T, \dots, (\mathbf{x}_N^*)^T]^T$ denote an optimal solution to (P2), and let $\boldsymbol{\nu}^* \triangleq \boldsymbol{\nu}_1^* = \dots = \boldsymbol{\nu}_N^*$ and $\{\mathbf{u}_{ij}^*, \mathbf{v}_{ij}^*\}$ denote a pair of optimal primal and dual solutions to problem (26) (i.e., (5)).

- (a) The sequence $\mathbf{x}^{(k)} = [(\mathbf{x}_1^{(k)})^T, \dots, (\mathbf{x}_N^{(k)})^T]^T$ generated from Algorithm 4 converges to \mathbf{x}^* of (P2) while $\boldsymbol{\nu}_1^{(k)}, \dots, \boldsymbol{\nu}_N^{(k)}$ converge to a common point $\boldsymbol{\nu}^*$ of problem (5).
- (b) If $\phi_i(\mathbf{x}) = f_i(\mathbf{A}_i \mathbf{x})$, where \mathbf{A}_i has full column rank, and \mathbf{E}_i has full row rank, for all $i \in V$, then for some $0 < \alpha < 1$ and $\rho > 0$, we have

$$\begin{aligned} & \|\mathbf{x}^{(k)} - \mathbf{x}^*\|_{\alpha \mathbf{M} + \frac{1}{2} \mathbf{P}}^2 + \frac{1}{c} \|\mathbf{u}^{(k+1)} - \mathbf{u}^*\|_2^2 \\ & + \frac{c}{2} \|\boldsymbol{\nu}^{(k)} - \mathbf{1}_N \otimes \boldsymbol{\nu}^*\|_{(\mathbf{D} + \mathbf{W}) \otimes \mathbf{I}_M}^2 \rightarrow 0 \text{ linearly,} \end{aligned} \quad (48)$$

Algorithm 4 Proposed IDC-ADMM for solving (P2)

1: **Given** initial variables $\mathbf{x}_i^{(0)} \in \mathbb{R}^K$ and $\mathbf{p}_i^{(0)} = \mathbf{0}$ for each agent i , $i \in V$. Set $k = 1$.

2: **repeat**

3: For all $i \in V$ (in parallel),

$$\begin{aligned} \mathbf{p}_i^{(k)} &= \mathbf{p}_i^{(k-1)} + c \sum_{j \in \mathcal{N}_i} (\boldsymbol{\nu}_i^{(k-1)} - \boldsymbol{\nu}_j^{(k-1)}), \\ \mathbf{x}_i^{(k)} &= \text{prox}_{g_i}^{\beta_i} \left[\mathbf{x}_i^{(k-1)} - \frac{1}{\beta_i} \mathbf{A}_i^T \nabla f_i(\mathbf{A}_i \mathbf{x}_i^{(k-1)}) \right. \\ &\quad \left. - \frac{1}{2\beta_i |\mathcal{N}_i|} \mathbf{E}_i^T \left(\frac{1}{c} (\mathbf{E}_i \mathbf{x}_i^{(k-1)} - \frac{1}{N} \mathbf{q}) - \frac{1}{c} \mathbf{p}_i^{(k)} \right) \right. \\ &\quad \left. + \sum_{j \in \mathcal{N}_i} (\boldsymbol{\nu}_i^{(k-1)} + \boldsymbol{\nu}_j^{(k-1)}) \right], \end{aligned} \quad (45)$$

$$\begin{aligned} \boldsymbol{\nu}_i^{(k)} &= \frac{1}{2|\mathcal{N}_i|} \left(\sum_{j \in \mathcal{N}_i} (\boldsymbol{\nu}_i^{(k-1)} + \boldsymbol{\nu}_j^{(k-1)}) - \frac{1}{c} \mathbf{p}_i^{(k)} \right) \\ &\quad + \frac{1}{c} (\mathbf{E}_i \mathbf{x}_i^{(k)} - \frac{1}{N} \mathbf{q}). \end{aligned} \quad (46)$$

4: **Set** $k = k + 1$.

5: **until** a predefined stopping criterion (e.g., a maximum iteration number) is satisfied.

where $\mathbf{u}^{(k)}$ and \mathbf{u}^* are defined similarly as in Theorem 1, \mathbf{M} is defined in (25), and $\mathbf{P} \triangleq \mathbf{D}_\beta - \frac{1}{2c} \text{blkdiag}\{\frac{1}{|\mathcal{N}_1|} \mathbf{E}_1^T \mathbf{E}_1, \dots, \frac{1}{|\mathcal{N}_N|} \mathbf{E}_N^T \mathbf{E}_N\} \succ \mathbf{0}$.

The proof is presented in Appendix B. Note that, in addition to the smooth and strongly convex objective function, IDC-ADMM also requires matrices \mathbf{E}_i 's to have full row rank in order to have a linear convergence rate.

V. NUMERICAL RESULTS

In this section, we examine the numerical performance of Algorithm 1 to 4 presented so far.

A. Performance of C-ADMM and IC-ADMM

To test C-ADMM (Algorithm 1) and IC-ADMM (Algorithm 2), we considered the distributed RPD LR problem in (7) with $\Psi_i(\mathbf{y}; \mathbf{A}_i, \mathbf{b}_i)$ in (8) and $g_i(\mathbf{y}) = \frac{\lambda}{N} \|\mathbf{y}\|_1 + \eta(\mathbf{y})$, where $\lambda > 0$ is a penalty parameter, and $\eta(\mathbf{y})$ is an indicator function specifying that the regression variables lie in a set $\mathcal{X} = \{\mathbf{y} \in \mathbb{R}^K \mid |x_i| \leq a \ \forall i\}$ for some $a > 0$ (see Eqn. (3)). We considered a simple two image classification task. Specifically, we used the images D24 and D68 from the Brodatz data set (<http://www.uix.no/~tranden/brodatz.html>) to generate the regression data matrix \mathbf{A} . We randomly extracted $(NM)/2$ overlapping patches with dimension $\sqrt{K} \times \sqrt{K}$ from the two images, respectively, followed by vectorizing the M patches into

vectors and stacking all of them into an $M \times K$ matrix. The rows of the matrix were randomly shuffled and the resultant matrix was used as the data matrix \mathbf{A} . For the RPD LR problem (7), we horizontally partitioned the matrix \mathbf{A} into N submatrices $\mathbf{A}_1, \dots, \mathbf{A}_N$, each with dimension $M \times K$. These matrices were used as the training data. Note that each \mathbf{A}_i contains patches from both images. The binary labels \mathbf{b}_i 's then were generated accordingly with 1 for one image and -1 for the other. The connected graph \mathcal{G} was randomly generated following the same method as in [39].

To implement C-ADMM (Algorithm 1), we employed the fast iterative shrinkage thresholding algorithm (FISTA) [40], [41] to solve subproblem (15) for each agent i . For (15), the associated FISTA steps can be shown as

$$\tilde{\mathbf{g}}_i^{(\ell)} = \max \left\{ -a, \min \left\{ a, \mathcal{S} \left[\mathbf{z}_i^{(\ell-1)} - \rho_i^{(\ell)} \left[\mathbf{A}_i^T \nabla f_i(\mathbf{A}_i \mathbf{z}_i^{(\ell-1)}) + \mathbf{p}_i^{(k)} + 2c \sum_{j \in \mathcal{N}_i} \left(\mathbf{z}_i^{(\ell-1)} - \frac{\mathbf{y}_i^{(k-1)} + \mathbf{y}_j^{(k-1)}}{2} \right) \right], \frac{\lambda \rho_i^{(\ell)}}{N} \right] \right\} \right\}, \quad (49a)$$

$$\mathbf{z}_i^{(\ell)} = \tilde{\mathbf{g}}_i^{(\ell)} + \frac{\ell-1}{\ell+2}(\tilde{\mathbf{g}}_i^{(\ell)} - \tilde{\mathbf{g}}_i^{(\ell-1)}), \quad (49b)$$

where ℓ denotes the inner iteration index of FISTA, $\rho_i^{(\ell)} > 0$ is a step size and \mathcal{S} is defined in (21). The stopping criterion of (49) was based on the PG residue (pgr) $\text{pgr} = \|\mathbf{z}_i^{(\ell-1)} - \tilde{\mathbf{g}}_i^{(\ell)}\| / (\rho_i^{(\ell)} \sqrt{K})$ [40], [41]. For obtaining a high-accuracy solution of (15), one may set the stopping criterion as, e.g., $\text{pgr} < 10^{-5}$. Suppose that FISTA stops at iteration $\ell_i(k)$. We then set $\mathbf{y}_i^{(k)} = \tilde{\mathbf{g}}_i^{(\ell_i(k))}$ as a solution to subproblem (15).

For IC-ADMM (Algorithm 2), the corresponding step in (20) is given by

$$\mathbf{y}_i^{(k)} = \max \left\{ -a, \min \left\{ a, \frac{1}{\gamma_i} \mathcal{S} \left[\beta \mathbf{y}_i^{(k-1)} - \mathbf{A}_i^T \nabla f_i(\mathbf{A}_i \mathbf{y}_i^{(k-1)}) - \mathbf{p}_i^{(k)} + c \sum_{j \in \mathcal{N}_i} (\mathbf{y}_i^{(k-1)} + \mathbf{y}_j^{(k-1)}), \frac{\lambda}{N} \right] \right\} \right\}. \quad (50)$$

From (??) and (49), the complexity of agent i at iteration k of C-ADMM is given by the order of $K + \ell_i(k)(2MK + 2K)$ if one counts only the multiplication operations; while from (??) and (50), the per-iteration complexity of each agent in IC-ADMM is given by the order of $K + (2MK + 2K)$. One can see that, for each agent i , the computational complexity of C-ADMM per iteration k (we refer this as the ‘‘ADMM iteration (ADMM Ite.)’’) is roughly $\ell_i(k)$ times that of IC-ADMM.

The stopping criterion of Algorithms 1 and 2 was based on measuring the solution accuracy $\text{acc} = (\text{obj}(\hat{\mathbf{y}}^{(k)}) - \text{obj}^*) / \text{obj}^*$ and variable consensus error $\text{cserr} = \sum_{i=1}^N \|\hat{\mathbf{y}}^{(k)} - \mathbf{y}_i^{(k)}\|_2^2 / N$, where $\hat{\mathbf{y}}^{(k)} = (\sum_{i=1}^N \mathbf{y}_i^{(k)}) / N$, $\text{obj}(\hat{\mathbf{y}}^{(k)})$ denotes the objective value of (7) given $\mathbf{y} = \hat{\mathbf{y}}^{(k)}$, and obj^* is the optimal

value of (7) which was obtained by FISTA [40], [41] with a high solution accuracy of $\text{pgr} < 10^{-6}$. The two algorithms were set to stop whenever acc and cserr are both smaller than preset target values.

In Table I(a), we considered a simulation example of $N = 10$, $K = 10,000$, $M = 10$, $\lambda = 0.1$ and $a = 1$, and display the comparison results. We not only present the required ADMM iterations but also the computation time per agent⁶ (in second) of the two methods. The convergence curves of C-ADMM and IC-ADMM with respect to the ADMM iteration are also shown in Figs. 1(a) and 1(b). The stopping conditions are $\text{acc} < 10^{-4}$ and $\text{cserr} < 10^{-5}$. For C-ADMM, we considered two cases, one with the stopping condition of FISTA for solving subproblem (15) set to $\text{pgr} < 10^{-5}$ and the other with that set to $\text{pgr} < 10^{-4}$. The penalty parameter c for C-ADMM was set to $c = 0.03$ and the step size $\rho_i^{(\ell)}$ of FISTA (see (49)) was set to a constant $\rho_i^{(\ell)} = 0.1$. The penalty parameters c and β of IC-ADMM were set to $c = 0.01$ and $\beta = 1.2$. We observe from Table I(a) that IC-ADMM in general requires more ADMM iterations than C-ADMM; however, the computation time is significantly smaller, as also illustrated in Figure 1(c). Specifically, the computation time of IC-ADMM is around $44.56/2.14 \approx 20.8$ times smaller than that of C-ADMM ($\text{pgr} < 10^{-5}$). We also observe that C-ADMM ($\text{pgr} < 10^{-4}$) consumes a smaller computation time for achieving $\text{acc} < 10^{-4}$. However, the associated $\text{cserr} = 3.425 \times 10^{-4}$ does not achieve the target value 10^{-5} . In fact, C-ADMM ($\text{pgr} < 10^{-4}$) cannot reduce cserr properly. As one can see from Fig. 1(b), the cserr curve of C-ADMM ($\text{pgr} < 10^{-4}$) keeps relatively high and does not decrease along the iterations. In Fig. 1(a) and Fig. 1(b), we also plot the convergence curves of the consensus subgradient method in [5], where the diminishing step size $10/k$ was used. As one can see, the consensus subgradient method converges much slower than IC-ADMM.

In Table I(b), we considered another example with the network size increased to $N = 50$. We set $c = 0.004$ for C-ADMM and $\rho_i^{(\ell)} = 0.1$ for FISTA; while for IC-ADMM, we set $c = 0.008$ and $\beta = 1.2$. The computation times of C-ADMM and IC-ADMM under this setting are also shown in Fig. 1(c). We can observe similar comparison results from Table I(b) and Fig. 1(c). Specifically, the computation time of IC-ADMM is around 8.75 times smaller than C-ADMM ($\text{pgr} < 10^{-5}$). When considering a lower accuracy of $\text{pgr} < 10^{-4}$, it is found that C-ADMM cannot properly converge.

To corroborating the linear convergence behavior of C-ADMM and IC-ADMM as claimed in Theorem 1(b)), we consider a problem instance of (7) with $\lambda = 0$, $N = 10$, $K = 25$, $M = 1,000$ and $a = 10$. We set $c = 0.2$ for C-ADMM and $\rho_i^{(\ell)} = 0.01$ and $\text{pgr} < 10^{-5}$ for FISTA; while for IC-ADMM, we set

⁶The simulation was performed on a desktop computer with 8-core Intel 1.3GHz CPU and 8 GB RAM. All the algorithms were implemented by MATLAB codes.

TABLE I: Comparison of C-ADMM and IC-ADMM

(a) $N = 10, K = 10,000, M = 10, \lambda = 0.1, a = 1.$			
	C-ADMM (pgr $< 10^{-5}$)	C-ADMM (pgr $< 10^{-4}$)	IC-ADMM
ADMM lte.	810	675	2973
Compt. Time (sec)	44.56	17.86	2.14
acc $< 10^{-4}$	9.982×10^{-5}	9.91×10^{-5}	9.99×10^{-5}
cserr $< 10^{-5}$	1.53×10^{-6}	3.425×10^{-4}	3.859×10^{-9}

(b) $N = 50, K = 10,000, M = 10, \lambda = 0.15, a = 1.$			
	C-ADMM (pgr $< 10^{-5}$)	C-ADMM (pgr $< 10^{-4}$)	IC-ADMM
ADMM lte.	952	N/A	7,251
Compt. Time (sec)	81.72	N/A	9.33
acc $< 10^{-4}$	9.99×10^{-5}	N/A	9.999×10^{-5}
cserr $< 10^{-5}$	1.305×10^{-7}	N/A	1.169×10^{-10}

$c = 1.2$ and $\beta = 10$. The convergence curves are shown in Figure 2. One can see from this figure that both algorithms converge linearly under this setting.

B. Performance of DC-ADMM and IDC-ADMM

We examine the performance of DC-ADMM (Algorithm 3) and IDC-ADMM (Algorithm 4) by considering the distributed CPD LR problem in (9), with $\Psi_i(\mathbf{x}_i; \mathbf{E}_i, \mathbf{b})$ in (10) and $g_i(\mathbf{x}_i) = \lambda \|\mathbf{x}_i\|_1$. Each variable \mathbf{x}_i is subject to the constraint set $\mathcal{X}_i = \{\mathbf{x}_i \in \mathbb{R}^{K/N} \mid |[\mathbf{x}_i]_j| \leq a \forall j\}$ for some $a > 0$. DC-ADMM and IDC-ADMM were applied to handle the associated problem (11). The regression data matrix $\mathbf{E} = [\mathbf{E}_1, \dots, \mathbf{E}_N]$ was generated following the same way as generating \mathbf{A} in Section V-A. To implement DC-ADMM, we employed FISTA [40], [41] to solve subproblem (35) and the solution accuracy was measured by the PG residue of FISTA.

In Table II(a), we show the comparison results for an example of $N = 50, K = 200, M = 100, \lambda = 0.05$ and $a = 10$. The convergence curves are also shown in Figs. 3(a) to 3(c). It was set $c = 0.05$ for DC-ADMM and the step size of FISTA $\rho_i^{(\ell)}$ was determined based on a line search rule [41]. We see from Table II(a) that, for achieving $\text{acc} < 10^{-4}$, DC-ADMM (pgr $< 10^{-5}$) took 329 ADMM iterations whereas IDC-ADMM took 10,814 iterations. However, the computation time of DC-ADMM (pgr $< 10^{-5}$)

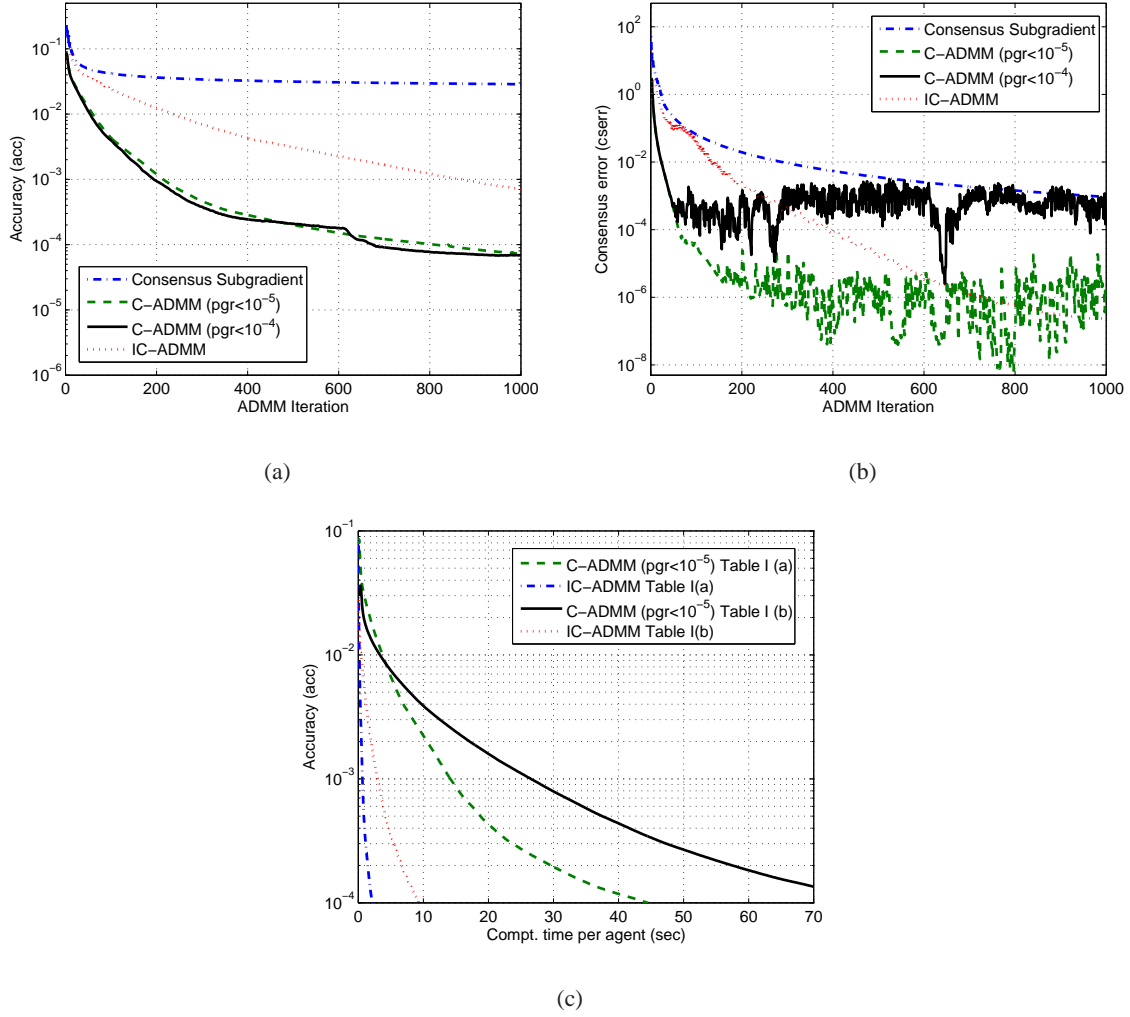


Fig. 1: Convergence curves of C-ADMM and IC-ADMM.

is around $42.78/1.92 \approx 22.28$ times higher than IDC-ADMM. When one reduce the solution accuracy of FISTA for solving subproblem (35) to $\text{pgr} < 10^{-4}$, DC-ADMM cannot reach the high accuracy of $\text{acc} < 10^{-4}$, as observed in Fig. 3(a). From Fig. 3(b), one can see that DC-ADMM converges much faster than IDC-ADMM with respect to the ADMM iterations. However, as shown from Fig. 3(c), the comparison result is reversed when one counts the computation times.

In Table II(b), we considered another example with K increased to 800. We set $c = 0.05$ for DC-ADMM, and set $c = 0.08$ and $\beta = 5$ for IDC-ADMM. From Table II(b) and Figs. 3(b) and 3(c), one can observe similar results.

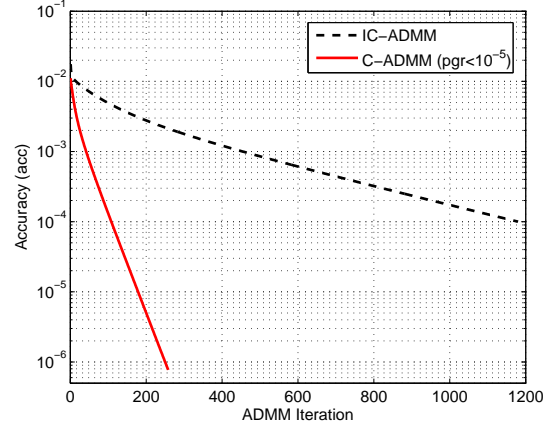


Fig. 2: Convergence curves of C-ADMM and IC-ADMM.

TABLE II: Comparison of DC-ADMM and IDC-ADMM

(a) $N = 50, K = 200, M = 100, \lambda = 0.05, a = 10$.

	DC-ADMM (pgr < 10^{-5})	DC-ADMM (pgr < 10^{-4})	IDC-ADMM
ADMM lte.	329	N/A	10814
Compt. Time (sec)	42.78	N/A	1.92
acc < 10^{-4}	9.928×10^{-5}	N/A	9.997×10^{-5}

(b) $N = 50, K = 800, M = 100, \lambda = 0.01, a = 20$.

	DC-ADMM (pgr < 10^{-5})	DC-ADMM (pgr < 10^{-4})	IDC-ADMM
ADMM lte.	475	N/A	38728
Compt. Time (sec)	427.73	N/A	18.07
acc < 10^{-4}	9.777×10^{-5}	N/A	9.999×10^{-5}

VI. CONCLUSIONS

In this paper, we have presented ADMM based distributed optimization methods for solving problems (P1) and (P2) in multi-agent networks. In particular, aiming at reducing the computational complexity of C-ADMM for solving large-scale instances of (P1) with complicated objective functions, we have proposed the IC-ADMM method (Algorithm 2) where agents perform one PG update only at each iteration. For (P2), we have proposed the DC-ADMM method (Algorithm 3) and its complexity reduced

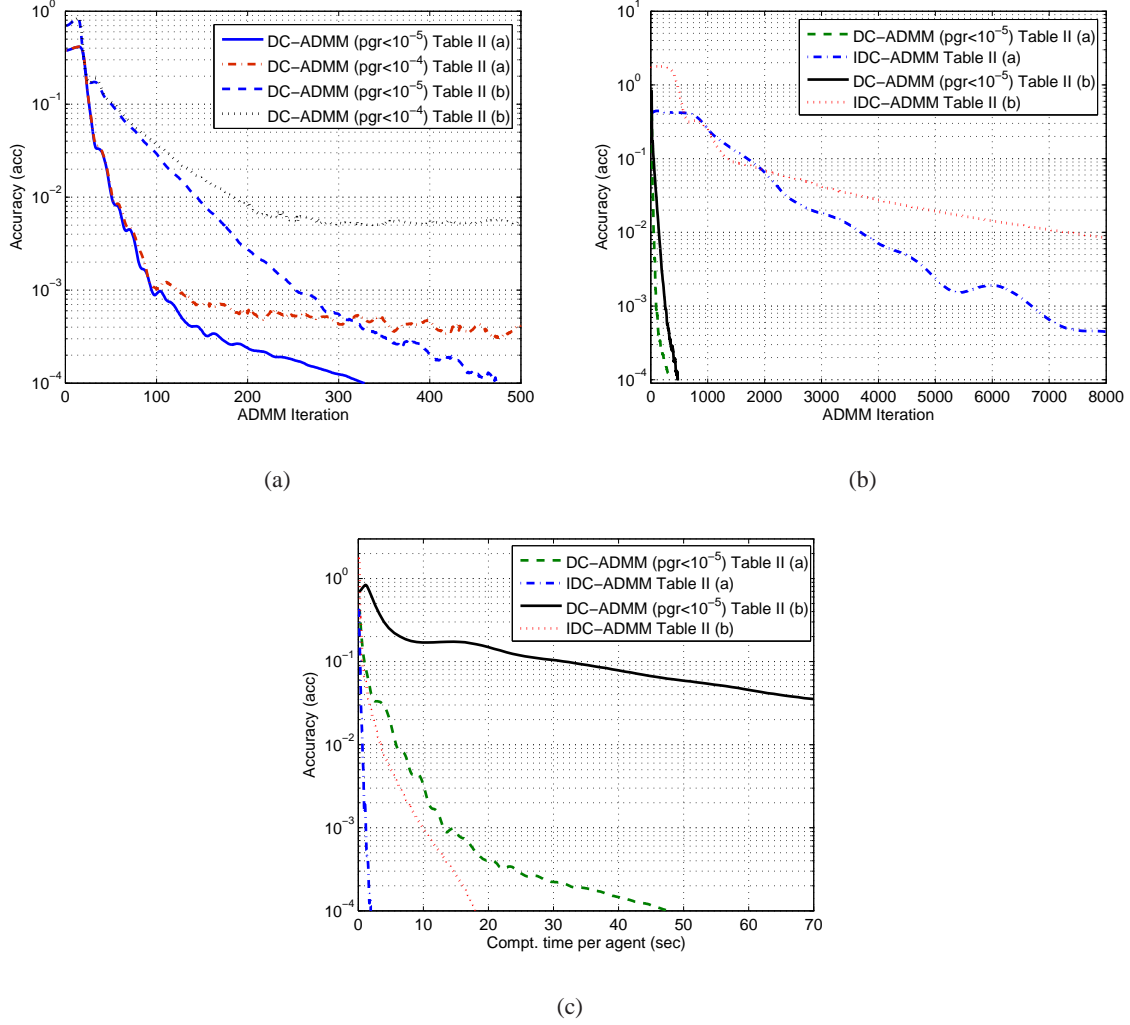


Fig. 3: Convergence curves of DC-ADMM and IDC-ADMM.

counterpart IDC-ADMM (Algorithm 4). Preliminary numerical results based on the distributed LR problems (7) and (11) have shown that the proposed methods converge faster than the consensus subgradient method. Moreover, both IC-ADMM and IDC-ADMM require more ADMM iterations than C-ADMM and DC-ADMM, but the traded computational complexity reduction is significant.

APPENDIX A

PROOF OF THEOREM 1

Proof of Theorem 1(a): Let $\tilde{\mathbf{y}}^* \triangleq [(\mathbf{y}_1^*)^T, \dots, (\mathbf{y}_N^*)^T]^T$ and $\{\mathbf{u}_{ij}^*, \mathbf{v}_{ij}^*, j \in \mathcal{N}_i\}_{i=1}^N$ be a pair of optimal primal and dual solutions to problem (13). Then they satisfy the following Karush-Kuhn-Tucker (KKT)

conditions: $\forall i \in V$,

$$\mathbf{A}_i^T \nabla f_i(\mathbf{A}_i \mathbf{y}_i^*) + \partial g_i(\mathbf{y}_i^*) + \sum_{j \in \mathcal{N}_i} (\mathbf{u}_{ij}^* + \mathbf{v}_{ji}^*) = \mathbf{0}, \quad (\text{A.1})$$

$$\mathbf{y}_i^* = \mathbf{y}_j^* \quad \forall j \in \mathcal{N}_i, \quad (\text{A.2})$$

$$\mathbf{u}_{ij}^* + \mathbf{v}_{ji}^* = \mathbf{0}, \quad \forall j \in \mathcal{N}_i, \quad (\text{A.3})$$

where $\partial g_i(\mathbf{y}_i^*)$ denotes the subgradient of g_i at \mathbf{y}_i^* . Under Assumption 1, (A.2) implies that $\mathbf{y}^* \triangleq \mathbf{y}_1^* = \dots = \mathbf{y}_N^*$ and $\tilde{\mathbf{y}}^* = \mathbf{1}_N \otimes \mathbf{y}^*$, i.e., consensus among agents is reached, and thus \mathbf{y}^* is optimal to the original problem (P1).

By recalling that $\mathbf{p}_i^{(k)} = \sum_{j \in \mathcal{N}_i} (\mathbf{u}_{ij}^{(k)} + \mathbf{v}_{ji}^{(k)}) \quad \forall i \in V$, and by the optimality condition of (18) [15], we have that

$$\begin{aligned} \mathbf{0} &= \mathbf{A}_i^T \nabla f_i(\mathbf{A}_i \mathbf{y}_i^{(k-1)}) + \beta_i(\mathbf{y}_i^{(k)} - \mathbf{y}_i^{(k-1)}) + \partial g_i(\mathbf{y}_i^{(k)}) \\ &\quad + \sum_{j \in \mathcal{N}_i} (\mathbf{u}_{ij}^{(k)} + \mathbf{v}_{ji}^{(k)}) \\ &\quad + 2c \sum_{j \in \mathcal{N}_i} \left(\mathbf{y}_i^{(k)} - \frac{\mathbf{y}_i^{(k-1)} + \mathbf{y}_j^{(k-1)}}{2} \right). \end{aligned} \quad (\text{A.4})$$

By combining (A.4) with (A.1), one obtains

$$\begin{aligned} &\mathbf{A}_i^T \nabla f_i(\mathbf{A}_i \mathbf{y}_i^{(k-1)}) - \mathbf{A}_i^T \nabla f_i(\mathbf{A}_i \mathbf{y}^*) + \beta_i(\mathbf{y}_i^{(k)} - \mathbf{y}_i^{(k-1)}) \\ &\quad + \partial g_i(\mathbf{y}_i^{(k)}) - \partial g_i(\mathbf{y}^*) + \sum_{j \in \mathcal{N}_i} (\mathbf{u}_{ij}^{(k)} + \mathbf{v}_{ji}^{(k)} - \mathbf{u}_{ij}^* - \mathbf{v}_{ji}^*) \\ &\quad + 2c \sum_{j \in \mathcal{N}_i} \left(\mathbf{y}_i^{(k)} - \frac{\mathbf{y}_i^{(k-1)} + \mathbf{y}_j^{(k-1)}}{2} \right) = \mathbf{0}. \end{aligned} \quad (\text{A.5})$$

Adding and subtracting $\mathbf{A}_i^T \nabla f_i(\mathbf{A}_i \mathbf{y}_i^{(k)})$ in the left hand side (LHS) of (A.5) followed by multiplying $(\mathbf{y}_i^{(k)} - \mathbf{y}^*)$ on both sides yields

$$\begin{aligned} &(\nabla f_i(\mathbf{A}_i \mathbf{y}_i^{(k-1)}) - \nabla f_i(\mathbf{A}_i \mathbf{y}_i^{(k)}))^T \mathbf{A}_i (\mathbf{y}_i^{(k)} - \mathbf{y}^*) + \beta_i(\mathbf{y}_i^{(k)} - \mathbf{y}_i^{(k-1)})^T (\mathbf{y}_i^{(k)} - \mathbf{y}^*) \\ &\quad + (\nabla f_i(\mathbf{A}_i \mathbf{y}_i^{(k)}) - \nabla f_i(\mathbf{A}_i \mathbf{y}^*))^T \mathbf{A}_i (\mathbf{y}_i^{(k)} - \mathbf{y}^*) \\ &\quad + (\partial g_i(\mathbf{y}_i^{(k)}) - \partial g_i(\mathbf{y}^*))^T (\mathbf{y}_i^{(k)} - \mathbf{y}^*) + \sum_{j \in \mathcal{N}_i} (\mathbf{u}_{ij}^{(k)} - \mathbf{u}_{ij}^* + \mathbf{v}_{ji}^{(k)} - \mathbf{v}_{ji}^*)^T (\mathbf{y}_i^{(k)} - \mathbf{y}^*) \\ &\quad + 2c \sum_{j \in \mathcal{N}_i} \left(\mathbf{y}_i^{(k)} - \frac{\mathbf{y}_i^{(k-1)} + \mathbf{y}_j^{(k-1)}}{2} \right)^T (\mathbf{y}_i^{(k)} - \mathbf{y}^*) = \mathbf{0}. \end{aligned} \quad (\text{A.6})$$

Note that the first term on the LHS of (A.6) can be lower bounded as

$$\begin{aligned}
& (\nabla f_i(\mathbf{A}_i \mathbf{y}_i^{(k-1)}) - \nabla f_i(\mathbf{A}_i \mathbf{y}_i^{(k)}))^T \mathbf{A}_i (\mathbf{y}_i^{(k)} - \mathbf{y}^*) \\
& \geq \frac{-1}{2\rho_i} \|\nabla f_i(\mathbf{A}_i \mathbf{y}_i^{(k-1)}) - \nabla f_i(\mathbf{A}_i \mathbf{y}_i^{(k)})\|_2^2 \\
& \quad - \frac{\rho_i}{2} \|\mathbf{y}_i^{(k)} - \mathbf{y}^*\|_{\mathbf{A}_i^T \mathbf{A}_i}^2 \\
& \geq \frac{-L_{f,i}^2}{2\rho_i} \|\mathbf{y}_i^{(k-1)} - \mathbf{y}_i^{(k)}\|_{\mathbf{A}_i^T \mathbf{A}_i}^2 - \frac{\rho_i}{2} \|\mathbf{y}_i^{(k)} - \mathbf{y}^*\|_{\mathbf{A}_i^T \mathbf{A}_i}^2
\end{aligned} \tag{A.7}$$

for any $\rho_i > 0$, where the second inequality is due to (12) in Assumption 3. By the strong convexity of f_i and convexity of g_i , the third and fourth terms of (A.6) can respectively be lower bounded as

$$\begin{aligned}
& (\nabla f_i(\mathbf{A}_i \mathbf{y}_i^{(k)}) - \nabla f_i(\mathbf{A}_i \mathbf{y}^*))^T \mathbf{A}_i (\mathbf{y}_i^{(k)} - \mathbf{y}^*) \\
& \geq \sigma_{f,i}^2 \|\mathbf{y}_i^{(k)} - \mathbf{y}^*\|_{\mathbf{A}_i^T \mathbf{A}_i}^2,
\end{aligned} \tag{A.8}$$

$$(\partial g_i(\mathbf{y}_i^{(k)}) - \partial g_i(\mathbf{y}^*))^T (\mathbf{y}_i^{(k)} - \mathbf{y}^*) \geq 0. \tag{A.9}$$

Moreover, it follows from (14a) and (14b) that the fifth term of (A.6) can be expressed as

$$\begin{aligned}
& \sum_{j \in \mathcal{N}_i} (\mathbf{u}_{ij}^{(k)} - \mathbf{u}_{ij}^* + \mathbf{v}_{ji}^{(k)} - \mathbf{v}_{ji}^*)^T (\mathbf{y}_i^{(k)} - \mathbf{y}^*) \\
& = \sum_{j \in \mathcal{N}_i} (\mathbf{u}_{ij}^{(k+1)} - \mathbf{u}_{ij}^* + \mathbf{v}_{ji}^{(k+1)} - \mathbf{v}_{ji}^*)^T (\mathbf{y}_i^{(k)} - \mathbf{y}^*) \\
& \quad - 2c \sum_{j \in \mathcal{N}_i} \left(\mathbf{y}_i^{(k)} - \frac{\mathbf{y}_i^{(k)} + \mathbf{y}_j^{(k)}}{2} \right)^T (\mathbf{y}_i^{(k)} - \mathbf{y}^*).
\end{aligned} \tag{A.10}$$

By substituting (A.7) to (A.10) into (A.6) and summing over $i = 1, \dots, N$, we obtain

$$\begin{aligned}
& \|\mathbf{y}^{(k)} - \tilde{\mathbf{y}}^*\|_{\mathbf{M}}^2 - \frac{1}{2} \|\mathbf{y}^{(k-1)} - \mathbf{y}^{(k)}\|_{\tilde{\mathbf{A}}^T \mathbf{D}_{L_f} \mathbf{D}_{\rho}^{-1} \tilde{\mathbf{A}}}^2 \\
& \quad + (\mathbf{y}^{(k)} - \mathbf{y}^{(k-1)})^T \mathbf{D}_{\beta} (\mathbf{y}^{(k)} - \tilde{\mathbf{y}}^*) \\
& + \sum_{i=1}^N \sum_{j \in \mathcal{N}_i} (\mathbf{u}_{ij}^{(k+1)} - \mathbf{u}_{ij}^*)^T (\mathbf{y}_i^{(k)} - \mathbf{y}^*) \\
& + \sum_{i=1}^N \sum_{j \in \mathcal{N}_i} (\mathbf{v}_{ji}^{(k+1)} - \mathbf{v}_{ji}^*)^T (\mathbf{y}_i^{(k)} - \mathbf{y}^*) \\
& + 2c \sum_{i=1}^N \sum_{j \in \mathcal{N}_i} \left(\frac{\mathbf{y}_i^{(k)} + \mathbf{y}_j^{(k)}}{2} - \frac{\mathbf{y}_i^{(k-1)} + \mathbf{y}_j^{(k-1)}}{2} \right)^T (\mathbf{y}_i^{(k)} - \mathbf{y}^*) \\
& \leq 0,
\end{aligned} \tag{A.11}$$

where $\mathbf{y}^{(k)} = [(\mathbf{y}_1^{(k)})^T, \dots, (\mathbf{y}_N^{(k)})^T]^T$, $\tilde{\mathbf{A}} = \text{blkdiag}\{\mathbf{A}_1, \dots, \mathbf{A}_N\}$, $\mathbf{D}_{L_f} = \text{diag}\{L_{f,1}^2, \dots, L_{f,N}^2\} \otimes \mathbf{I}_K$, $\mathbf{D}_\beta = \text{diag}\{\beta_1, \dots, \beta_N\} \otimes \mathbf{I}_K$, $\mathbf{D}_\rho = \text{diag}\{\rho_1, \dots, \rho_N\} \otimes \mathbf{I}_K$, and as defined in (25),

$$\mathbf{M} = \tilde{\mathbf{A}}^T (\mathbf{D}_{\sigma_f} - \frac{1}{2} \mathbf{D}_\rho) \tilde{\mathbf{A}}.$$

It can be observed from (A.3) and also (14a) and (14b) that

$$\mathbf{u}_{ij}^* + \mathbf{v}_{ij}^* = \mathbf{0} \quad \forall j, i, \quad (\text{A.12})$$

$$\mathbf{u}_{ij}^{(k)} + \mathbf{v}_{ij}^{(k)} = \mathbf{0} \quad \forall j, i, k, \quad (\text{A.13})$$

given the initial $\mathbf{u}_{ij}^{(0)} + \mathbf{v}_{ij}^{(0)} = \mathbf{0} \quad \forall j, i, k$ which is equivalent to setting $\mathbf{p}_i^{(k)} = \mathbf{0} \quad \forall i \in V$ (See Step 1 of Algorithm 2). Besides, due to the symmetric property of \mathbf{W} , for any $\{\alpha_{ij}\}$, we have

$$\begin{aligned} \sum_{i=1}^N \sum_{j \in \mathcal{N}_i} \alpha_{ij} &= \sum_{i=1}^N \sum_{j=1}^N [\mathbf{W}]_{i,j} \alpha_{ij} \\ &= \sum_{i=1}^N \sum_{j=1}^N [\mathbf{W}]_{i,j} \alpha_{ji} = \sum_{i=1}^N \sum_{j \in \mathcal{N}_i} \alpha_{ji}. \end{aligned} \quad (\text{A.14})$$

By the above two properties, the fourth and fifth terms in the LHS of (A.11) can be written as

$$\begin{aligned} &\sum_{i=1}^N \sum_{j \in \mathcal{N}_i} (\mathbf{u}_{ij}^{(k+1)} - \mathbf{u}_{ij}^*)^T (\mathbf{y}_i^{(k)} - \mathbf{y}^*) \\ &\quad + \sum_{i=1}^N \sum_{j \in \mathcal{N}_i} (\mathbf{v}_{ji}^{(k+1)} - \mathbf{v}_{ji}^*)^T (\mathbf{y}_i^{(k)} - \mathbf{y}^*) \\ &= \sum_{i=1}^N \sum_{j \in \mathcal{N}_i} (\mathbf{u}_{ij}^{(k+1)} - \mathbf{u}_{ij}^*)^T (\mathbf{y}_i^{(k)} - \mathbf{y}^*) \\ &\quad + \sum_{i=1}^N \sum_{j \in \mathcal{N}_i} (\mathbf{v}_{ij}^{(k+1)} - \mathbf{v}_{ij}^*)^T (\mathbf{y}_j^{(k)} - \mathbf{y}^*) \\ &= \sum_{i=1}^N \sum_{j \in \mathcal{N}_i} (\mathbf{u}_{ij}^{(k+1)} - \mathbf{u}_{ij}^*)^T (\mathbf{y}_i^{(k)} - \mathbf{y}_j^{(k)}) \\ &= \frac{2}{c} \sum_{i=1}^N \sum_{j \in \mathcal{N}_i} (\mathbf{u}_{ij}^{(k+1)} - \mathbf{u}_{ij}^*)^T (\mathbf{u}_{ij}^{(k+1)} - \mathbf{u}_{ij}^{(k)}) \\ &\triangleq \frac{2}{c} (\mathbf{u}^{(k+1)} - \mathbf{u}^*)^T (\mathbf{u}^{(k+1)} - \mathbf{u}^{(k)}), \end{aligned} \quad (\text{A.15})$$

where the first equality is owing to (A.14), the second equality is by (A.12) and (A.13), and the third equality is due to (14a). In (A.15), $\mathbf{u}^{(k)}$ (\mathbf{u}^*) is a vector that stacks $\mathbf{u}_{ij}^{(k)}$ (\mathbf{u}_{ij}^*) for all $j \in \mathcal{N}_i$, $i = 1, \dots, N$.

The sixth term in the LHS of (A.11) can be rearranged as follows

$$\begin{aligned}
& c \sum_{i=1}^N \sum_{j \in \mathcal{N}_i} (\mathbf{y}_i^{(k)} - \mathbf{y}_i^{(k-1)})^T (\mathbf{y}_i^{(k)} - \mathbf{y}^*) \\
& \quad + c \sum_{i=1}^N \sum_{j \in \mathcal{N}_i} (\mathbf{y}_j^{(k)} - \mathbf{y}_j^{(k-1)})^T (\mathbf{y}_i^{(k)} - \mathbf{y}^*) \\
& = c \sum_{i=1}^N |\mathcal{N}_i| (\mathbf{y}_i^{(k)} - \mathbf{y}_i^{(k-1)})^T (\mathbf{y}_i^{(k)} - \mathbf{y}^*) \\
& \quad + c \sum_{i=1}^N \sum_{j=1}^N [\mathbf{W}]_{i,j} (\mathbf{y}_j^{(k)} - \mathbf{y}_j^{(k-1)})^T (\mathbf{y}_i^{(k)} - \mathbf{y}^*) \\
& = c (\mathbf{y}^{(k)} - \mathbf{y}^{(k-1)})^T (\mathbf{D} \otimes \mathbf{I}_K) (\mathbf{y}^{(k)} - \tilde{\mathbf{y}}^*) \\
& \quad + c (\mathbf{y}^{(k)} - \mathbf{y}^{(k-1)})^T (\mathbf{W} \otimes \mathbf{I}_K) (\mathbf{y}^{(k)} - \tilde{\mathbf{y}}^*) \\
& = c (\mathbf{y}^{(k)} - \mathbf{y}^{(k-1)})^T [(\mathbf{D} + \mathbf{W}) \otimes \mathbf{I}_K] (\mathbf{y}^{(k)} - \tilde{\mathbf{y}}^*). \tag{A.16}
\end{aligned}$$

Note that, by the graph theory [32], the normalized Laplacian matrix, i.e., $\mathbf{D}^{-\frac{1}{2}} \mathbf{L} \mathbf{D}^{-\frac{1}{2}}$, have $\lambda_{\max}(\mathbf{D}^{-\frac{1}{2}} \mathbf{L} \mathbf{D}^{-\frac{1}{2}}) \leq 2$. Thus, in (A.16),

$$\mathbf{D} + \mathbf{W} = 2\mathbf{D} - \mathbf{L} = \mathbf{D}^{\frac{1}{2}} (2\mathbf{I}_N - \mathbf{D}^{-\frac{1}{2}} \mathbf{L} \mathbf{D}^{-\frac{1}{2}}) \mathbf{D}^{\frac{1}{2}} \succeq \mathbf{0}.$$

By substituting (A.15) and (A.16) into (A.11), we obtain

$$\begin{aligned}
& \|\mathbf{y}^{(k)} - \tilde{\mathbf{y}}^*\|_{\mathbf{M}}^2 - \frac{1}{2} \|\mathbf{y}^{(k-1)} - \mathbf{y}^{(k)}\|_{\tilde{\mathbf{A}}^T \mathbf{D}_{L_f} \mathbf{D}_{\rho}^{-1} \tilde{\mathbf{A}}}^2 \\
& \quad + (\mathbf{y}^{(k)} - \mathbf{y}^{(k-1)})^T \mathbf{G} (\mathbf{y}^{(k)} - \tilde{\mathbf{y}}^*) \\
& \quad + \frac{2}{c} (\mathbf{u}^{(k+1)} - \mathbf{u}^*)^T (\mathbf{u}^{(k+1)} - \mathbf{u}^{(k)}) \leq 0, \tag{A.17}
\end{aligned}$$

where as defined in (24),

$$\mathbf{G} \triangleq \mathbf{D}_{\beta} + c((\mathbf{D} + \mathbf{W}) \otimes \mathbf{I}_K) \succ \mathbf{0}.$$

Note that

$$\begin{aligned}
& (\mathbf{a}^{(k)} - \mathbf{a}^{(k-1)})^T \mathbf{Q} (\mathbf{a}^{(k)} - \mathbf{a}^*) = \frac{1}{2} \|\mathbf{a}^{(k)} - \mathbf{a}^*\|_{\mathbf{Q}}^2 \\
& \quad + \frac{1}{2} \|\mathbf{a}^{(k)} - \mathbf{a}^{(k-1)}\|_{\mathbf{Q}}^2 - \frac{1}{2} \|\mathbf{a}^{(k-1)} - \mathbf{a}^*\|_{\mathbf{Q}}^2 \tag{A.18}
\end{aligned}$$

for any sequence $\mathbf{a}^{(k)}$ and matrix $\mathbf{Q} \succeq \mathbf{0}$. By applying (A.18) to each of the terms in (A.17), one obtains that

$$\begin{aligned} & (\mathbf{y}^{(k)} - \tilde{\mathbf{y}}^*)^T \left[\mathbf{M} + \frac{1}{2} \mathbf{G} \right] (\mathbf{y}^{(k)} - \tilde{\mathbf{y}}^*) + \frac{1}{c} \|\mathbf{u}^{(k+1)} - \mathbf{u}^*\|_2^2 \\ & \leq \frac{1}{2} (\mathbf{y}^{(k-1)} - \tilde{\mathbf{y}}^*)^T \mathbf{G} (\mathbf{y}^{(k-1)} - \tilde{\mathbf{y}}^*) \\ & \quad + \frac{1}{c} \|\mathbf{u}^{(k)} - \mathbf{u}^*\|_2^2 - \frac{1}{c} \|\mathbf{u}^{(k+1)} - \mathbf{u}^{(k)}\|_2^2 \\ & \quad - (\mathbf{y}^{(k)} - \mathbf{y}^{(k-1)})^T \left[\frac{1}{2} \mathbf{G} - \frac{1}{2} \tilde{\mathbf{A}}^T \mathbf{D}_{L_f} \mathbf{D}_\rho^{-1} \tilde{\mathbf{A}} \right] (\mathbf{y}^{(k)} - \mathbf{y}^{(k-1)}). \end{aligned} \quad (\text{A.19})$$

Now, consider the condition on β_i in (23). It can be easily checked that (23) implies that

$$\sigma_{f,i}^2 - \frac{\rho_i}{2} > 0, \quad (\text{A.20a})$$

$$\beta_i \mathbf{I}_K + c \lambda_{\min}(\mathbf{D} + \mathbf{W}) \mathbf{I}_K - \frac{L_{f,i}^2}{\rho_i} \mathbf{A}_i^T \mathbf{A}_i \succ \mathbf{0}, \quad (\text{A.20b})$$

for some $\sigma_{f,i}^2 \leq \rho_i < 2\sigma_{f,i}^2 \forall i \in V$, and therefore

$$\mathbf{M} \succeq \mathbf{0}, \quad \mathbf{G} - \tilde{\mathbf{A}}^T \mathbf{D}_{L_f} \mathbf{D}_\rho^{-1} \tilde{\mathbf{A}} \succ \mathbf{0}. \quad (\text{A.21})$$

With (A.21), (A.19) implies the following two results (R1) as $k \rightarrow \infty$, the sequence $\frac{1}{2} \|\mathbf{y}^{(k)} - \tilde{\mathbf{y}}^*\|_G^2 + \frac{1}{c} \|\mathbf{u}^{(k+1)} - \mathbf{u}^*\|_2^2$ converges for any pair of optimal $\tilde{\mathbf{y}}^*$ and \mathbf{u}^* to problem (13); and (R2)

$$\mathbf{y}^{(k)} - \mathbf{y}^{(k-1)} \rightarrow \mathbf{0}, \quad \mathbf{u}^{(k+1)} - \mathbf{u}^{(k)} \rightarrow \mathbf{0}. \quad (\text{A.22})$$

The result (R1) implies that the sequences of $\{\mathbf{y}_i^{(k)}\}$ and $\{\mathbf{u}_{ij}^{(k)}\}$ (so is $\{\mathbf{v}_{ij}^{(k)}\}$) are bounded. Let $\tilde{\mathbf{y}} = [(\hat{\mathbf{y}}_1)^T, \dots, (\hat{\mathbf{y}}_N)^T]^T$, $\hat{\mathbf{u}}_{ij}$ and $\hat{\mathbf{v}}_{ij}$ be a set of limit points of $\{\mathbf{y}^{(k)}\}$, $\{\mathbf{u}_{ij}^{(k)}\}$ and $\{\mathbf{v}_{ij}^{(k)}\}$, respectively. Firstly, by the result of $\mathbf{u}^{(k+1)} - \mathbf{u}^{(k)} \rightarrow \mathbf{0}$ and (14a), we have

$$\mathbf{y}_i^{(k)} - \mathbf{y}_j^{(k)} \rightarrow \mathbf{0} \implies \hat{\mathbf{y}} \triangleq \hat{\mathbf{y}}_i = \hat{\mathbf{y}}_j, \quad \forall j, i. \quad (\text{A.23})$$

Secondly, by (A.13), we have

$$\hat{\mathbf{u}}_{ij} + \hat{\mathbf{v}}_{ij} = \mathbf{0} \quad \forall j, i. \quad (\text{A.24})$$

Thirdly, by applying the result of $\mathbf{y}^{(k)} - \mathbf{y}^{(k-1)} \rightarrow \mathbf{0}$ and (A.23) to (A.4), we have

$$\mathbf{0} = \mathbf{A}_i^T \nabla f_i(\mathbf{A}_i \hat{\mathbf{y}}_i) + \partial g_i(\hat{\mathbf{y}}_i) + \sum_{j \in \mathcal{N}_i} (\hat{\mathbf{u}}_{ij} + \hat{\mathbf{v}}_{ji}) \quad (\text{A.25})$$

for all $i \in V$. So, $\tilde{\mathbf{y}}$ and $\{\hat{\mathbf{u}}_{ij}, \hat{\mathbf{v}}_{ij}\}$ are in fact a pair of optimal primal and dual solutions to problem (13) [see (A.1), (A.2) and (A.3)]. Therefore, according to (R1), the sequence $\frac{1}{2} \|\mathbf{y}^{(k)} - \tilde{\mathbf{y}}\|_G^2 + \frac{1}{c} \|\mathbf{u}^{(k+1)} - \hat{\mathbf{u}}\|_2^2$ converges. Furthermore, since $\frac{1}{2} \|\mathbf{y}^{(k)} - \tilde{\mathbf{y}}\|_G^2 + \frac{1}{c} \|\mathbf{u}^{(k+1)} - \hat{\mathbf{u}}\|_2^2$ has a limit value equal to zero, we conclude

that $\frac{1}{2}\|\mathbf{y}^{(k)} - \tilde{\mathbf{y}}\|_{\mathbf{G}}^2 + \frac{1}{c}\|\mathbf{u}^{(k+1)} - \hat{\mathbf{u}}\|_2^2$ in fact converges to zero. This says that $\mathbf{y}_i^{(k)} \rightarrow \hat{\mathbf{y}} \forall i \in V$ and $\mathbf{u}^{(k+1)} \rightarrow \hat{\mathbf{u}}$. The proof is thus complete. \blacksquare

Proof of Theorem 1(b): Let $0 < \alpha < 1$ be some positive number and rewrite (A.19) as

$$\begin{aligned} & \left(\|\mathbf{y}^{(k)} - \tilde{\mathbf{y}}^*\|_{\frac{1}{2}\mathbf{G}+\alpha\mathbf{M}}^2 + \frac{1}{c}\|\mathbf{u}^{(k+1)} - \mathbf{u}^*\|_2^2 \right) + \|\mathbf{y}^{(k)} - \tilde{\mathbf{y}}^*\|_{(1-\alpha)\mathbf{M}}^2 \\ & + \|\mathbf{y}^{(k-1)} - \tilde{\mathbf{y}}^*\|_{\alpha\mathbf{M}}^2 + \frac{1}{c}\|\mathbf{u}^{(k+1)} - \mathbf{u}^{(k)}\|_2^2 \\ & + (\mathbf{y}^{(k)} - \mathbf{y}^{(k-1)})^T \left[\frac{1}{2}\mathbf{G} - \frac{1}{2}\tilde{\mathbf{A}}^T \mathbf{D}_{L_f} \mathbf{D}_\rho^{-1} \tilde{\mathbf{A}} \right] (\mathbf{y}^{(k)} - \mathbf{y}^{(k-1)}) \\ & \leq \left(\|\mathbf{y}^{(k-1)} - \tilde{\mathbf{y}}^*\|_{\frac{1}{2}\mathbf{G}+\alpha\mathbf{M}}^2 + \frac{1}{c}\|\mathbf{u}^{(k)} - \mathbf{u}^*\|_2^2 \right). \end{aligned} \quad (\text{A.26})$$

Then, in order to prove linear convergence rate, i.e., for some $\delta > 0$,

$$\begin{aligned} & \left(\|\mathbf{y}^{(k)} - \tilde{\mathbf{y}}^*\|_{\frac{1}{2}\mathbf{G}+\alpha\mathbf{M}}^2 + \frac{1}{c}\|\mathbf{u}^{(k+1)} - \mathbf{u}^*\|_2^2 \right) \\ & \leq \frac{1}{1+\delta} \left(\|\mathbf{y}^{(k-1)} - \tilde{\mathbf{y}}^*\|_{\frac{1}{2}\mathbf{G}+\alpha\mathbf{M}}^2 + \frac{1}{c}\|\mathbf{u}^{(k)} - \mathbf{u}^*\|_2^2 \right), \end{aligned}$$

it is sufficient to show that

$$\begin{aligned} & \|\mathbf{y}^{(k)} - \tilde{\mathbf{y}}^*\|_{(1-\alpha)\mathbf{M}}^2 + \|\mathbf{y}^{(k-1)} - \tilde{\mathbf{y}}^*\|_{\alpha\mathbf{M}}^2 + \frac{1}{c}\|\mathbf{u}^{(k+1)} - \mathbf{u}^{(k)}\|_2^2 \\ & + (\mathbf{y}^{(k)} - \mathbf{y}^{(k-1)})^T \left[\frac{1}{2}\mathbf{G} - \frac{1}{2}\tilde{\mathbf{A}}^T \mathbf{D}_{L_f} \mathbf{D}_\rho^{-1} \tilde{\mathbf{A}} \right] (\mathbf{y}^{(k)} - \mathbf{y}^{(k-1)}) \\ & \geq \delta \left(\|\mathbf{y}^{(k)} - \tilde{\mathbf{y}}^*\|_{\frac{1}{2}\mathbf{G}+\alpha\mathbf{M}}^2 + \frac{1}{c}\|\mathbf{u}^{(k+1)} - \mathbf{u}^*\|_2^2 \right). \end{aligned} \quad (\text{A.27})$$

Recall from (A.5) and (A.10) that

$$\begin{aligned} & \mathbf{A}_i^T \nabla f_i(\mathbf{A}_i \mathbf{y}_i^{(k-1)}) - \mathbf{A}_i^T \nabla f_i(\mathbf{A}_i \mathbf{y}^*) + \beta_i (\mathbf{y}_i^{(k)} - \mathbf{y}_i^{(k-1)}) \\ & + \sum_{j \in \mathcal{N}_i} (\mathbf{u}_{ij}^{(k+1)} - \mathbf{u}_{ij}^*) + \sum_{j \in \mathcal{N}_i} (\mathbf{v}_{ji}^{(k+1)} - \mathbf{v}_{ji}^*) \\ & + 2c \sum_{j \in \mathcal{N}_i} \left(\frac{\mathbf{y}_i^{(k)} + \mathbf{y}_j^{(k)}}{2} - \frac{\mathbf{y}_i^{(k-1)} + \mathbf{y}_j^{(k-1)}}{2} \right) = \mathbf{0}. \end{aligned} \quad (\text{A.28})$$

By applying (A.12) and (A.13), (A.28) can be expressed as

$$\begin{aligned} & \mathbf{A}_i^T \nabla f_i(\mathbf{A}_i \mathbf{y}_i^{(k-1)}) - \mathbf{A}_i^T \nabla f_i(\mathbf{A}_i \mathbf{y}^*) + \beta_i (\mathbf{y}_i^{(k)} - \mathbf{y}_i^{(k-1)}) \\ & + \sum_{j \in \mathcal{N}_i} (\mathbf{u}_{ij}^{(k+1)} - \mathbf{u}_{ij}^{(k)} - \mathbf{u}_{ij}^* + \mathbf{u}_{ji}^*) \\ & + c \sum_{j \in \mathcal{N}_i} \left(\mathbf{y}_i^{(k)} - \mathbf{y}_i^{(k-1)} + \mathbf{y}_j^{(k)} - \mathbf{y}_j^{(k-1)} \right) = \mathbf{0}. \end{aligned} \quad (\text{A.29})$$

After stacking (A.29) for $i = 1, \dots, N$, one obtains

$$\begin{aligned} & \tilde{\mathbf{A}}^T (\nabla \mathbf{f}(\tilde{\mathbf{A}} \mathbf{y}^{(k-1)}) - \nabla \mathbf{f}(\tilde{\mathbf{A}} \mathbf{y}^*)) + \mathbf{G} (\mathbf{y}^{(k)} - \mathbf{y}^{(k-1)}) \\ & + \mathbf{\Upsilon} (\mathbf{u}^{(k+1)} - \mathbf{u}^*) = \mathbf{0}. \end{aligned} \quad (\text{A.30})$$

where $\nabla \mathbf{f}(\tilde{\mathbf{A}}\mathbf{y}^{(k)}) \triangleq [(\nabla f_1(\mathbf{A}_1\mathbf{y}_1^{(k)}))^T, \dots, (\nabla f_1(\mathbf{A}_N\mathbf{y}_N^{(k)}))^T]^T$ and $\mathbf{\Upsilon} \in \mathbb{R}^{KN \times 2|\mathcal{E}|K}$ is a linear mapping matrix satisfying

$$\begin{bmatrix} \sum_{j \in \mathcal{N}_1} (\mathbf{u}_{1j}^{(k+1)} - \mathbf{u}_{j1}^{(k+1)}) \\ \vdots \\ \sum_{j \in \mathcal{N}_N} (\mathbf{u}_{Nj}^{(k+1)} - \mathbf{u}_{jN}^{(k+1)}) \end{bmatrix} = \mathbf{\Upsilon} \mathbf{u}^{(k+1)}. \quad (\text{A.31})$$

According to [26]⁷, both $\mathbf{u}^{(k+1)}$ and \mathbf{u}^* lie in the range space of $\mathbf{\Upsilon}^T$. Hence, one can show that

$$\|\mathbf{\Upsilon}(\mathbf{u}^{(k+1)} - \mathbf{u}^*)\|^2 \geq \sigma_{\min}^2(\mathbf{\Upsilon}) \|\mathbf{u}^{(k+1)} - \mathbf{u}^*\|_2^2 \quad (\text{A.32})$$

where $\sigma_{\min}(\mathbf{\Upsilon}) > 0$ is the minimum nonzero singular value of $\mathbf{\Upsilon}$. From (A.30), we have that

$$\begin{aligned} & \|\mathbf{G}(\mathbf{y}^{(k)} - \mathbf{y}^{(k-1)})\|_2^2 \\ &= \|\mathbf{A}^T(\nabla \mathbf{f}(\tilde{\mathbf{A}}\mathbf{y}^{(k-1)}) - \nabla \mathbf{f}(\tilde{\mathbf{A}}\tilde{\mathbf{y}}^*)) - \mathbf{\Upsilon}(\mathbf{u}^{(k+1)} - \mathbf{u}^*)\|_2^2 \\ &\geq (1 - \mu) \|\mathbf{A}^T(\nabla \mathbf{f}(\tilde{\mathbf{A}}\mathbf{y}^{(k)}) - \nabla \mathbf{f}(\tilde{\mathbf{A}}\tilde{\mathbf{y}}^*))\|_2^2 \\ &\quad + (1 - \frac{1}{\mu}) \|\mathbf{\Upsilon}(\mathbf{u}^{(k+1)} - \mathbf{u}^*)\|_2^2 \\ &\geq (1 - \mu) \lambda_{\max}(\tilde{\mathbf{A}}^T \tilde{\mathbf{A}}) \|\nabla \mathbf{f}(\tilde{\mathbf{A}}\mathbf{y}^{(k)}) - \nabla \mathbf{f}(\tilde{\mathbf{A}}\tilde{\mathbf{y}}^*)\|_2^2 \\ &\quad + (1 - \frac{1}{\mu}) \sigma_{\min}^2(\mathbf{\Upsilon}) \|\mathbf{u}^{(k+1)} - \mathbf{u}^*\|_2^2 \\ &\geq (1 - \mu) \lambda_{\max}(\tilde{\mathbf{A}}^T \tilde{\mathbf{A}}) \|\mathbf{y}^{(k-1)} - \mathbf{y}^*\|_{\tilde{\mathbf{A}}^T \mathbf{D}_{L_f} \tilde{\mathbf{A}}}^2 \\ &\quad + (1 - \frac{1}{\mu}) \sigma_{\min}^2(\mathbf{\Upsilon}) \|\mathbf{u}^{(k+1)} - \mathbf{u}^*\|_2^2, \end{aligned} \quad (\text{A.33})$$

where the first inequality is due to the fact that

$$\|\mathbf{a} + \mathbf{q}\|_2^2 \geq (1 - \mu) \|\mathbf{a}\|_2^2 + (1 - \frac{1}{\mu}) \|\mathbf{q}\|_2^2 \quad (\text{A.34})$$

for any \mathbf{a}, \mathbf{q} and $\mu > 0$, the second inequality is obtained by setting $\mu > 1$ and (A.32), and the last inequality is by (12). Equation (A.33) implies that

$$\begin{aligned} & \frac{\delta}{c} \|\mathbf{u}^{(k+1)} - \mathbf{u}^*\|_2^2 \leq \frac{\delta}{c(1 - \frac{1}{\mu}) \sigma_{\min}^2(\mathbf{\Upsilon})} \|\mathbf{y}^{(k)} - \mathbf{y}^{(k-1)}\|_{\mathbf{G}^T \mathbf{G}}^2 \\ & + \frac{\delta(\mu - 1) \lambda_{\max}(\tilde{\mathbf{A}}^T \tilde{\mathbf{A}})}{c(1 - \frac{1}{\mu}) \sigma_{\min}^2(\mathbf{\Upsilon})} \|\mathbf{y}^{(k-1)} - \tilde{\mathbf{y}}^*\|_{\tilde{\mathbf{A}}^T \mathbf{D}_{L_f} \tilde{\mathbf{A}}}^2. \end{aligned} \quad (\text{A.35})$$

⁷Note that the matrix $\mathbf{\Upsilon}$ corresponds to matrix M_- in [26].

According to (A.35), (A.27) can hold true if

$$\begin{aligned}
& \|\mathbf{y}^{(k)} - \tilde{\mathbf{y}}^*\|_{(1-\alpha)\mathbf{M}}^2 \geq \delta \|\mathbf{y}^{(k)} - \tilde{\mathbf{y}}^*\|_{\frac{1}{2}\mathbf{G} + \alpha\mathbf{M}}^2, \\
& (\mathbf{y}^{(k)} - \mathbf{y}^{(k-1)})^T \left[\frac{1}{2}\mathbf{G} - \frac{1}{2}\tilde{\mathbf{A}}^T \mathbf{D}_{L_f} \mathbf{D}_\rho^{-1} \tilde{\mathbf{A}} \right] (\mathbf{y}^{(k)} - \mathbf{y}^{(k-1)}) \\
& \geq \frac{\delta}{c(1 - \frac{1}{\mu})\sigma_{\min}^2(\Upsilon)} \|\mathbf{y}^{(k)} - \mathbf{y}^{(k-1)}\|_{\mathbf{G}^T \mathbf{G}}^2, \\
& \|\mathbf{y}^{(k-1)} - \tilde{\mathbf{y}}^*\|_{\alpha\mathbf{M}}^2 \\
& \geq \frac{\delta(\mu - 1)\lambda_{\max}(\tilde{\mathbf{A}}^T \tilde{\mathbf{A}})}{c(1 - \frac{1}{\mu})\sigma_{\min}^2(\Upsilon)} \|(\mathbf{y}^{(k-1)} - \tilde{\mathbf{y}}^*)\|_{\tilde{\mathbf{A}}^T \mathbf{D}_{L_f} \tilde{\mathbf{A}}}^2,
\end{aligned}$$

which are respectively satisfied if the following three conditions can be satisfied for some $\delta > 0$

$$(1 - \alpha)\mathbf{M} \succeq \delta \left(\frac{1}{2}\mathbf{G} + \alpha\mathbf{M} \right), \quad (\text{A.36a})$$

$$\frac{1}{2}\mathbf{G} - \frac{1}{2}\tilde{\mathbf{A}}^T \mathbf{D}_{L_f} \mathbf{D}_\rho^{-1} \tilde{\mathbf{A}} \succeq \frac{\delta}{c(1 - \frac{1}{\mu})\sigma_{\min}^2(\Upsilon)} \mathbf{G}^T \mathbf{G}, \quad (\text{A.36b})$$

$$\alpha(\mathbf{D}_{\sigma_f} - \frac{1}{2}\mathbf{D}_\rho) \succeq \delta \frac{(\mu - 1)\lambda_{\max}(\tilde{\mathbf{A}}^T \tilde{\mathbf{A}})}{c(1 - \frac{1}{\mu})\sigma_{\min}^2(\Upsilon)} \mathbf{D}_{L_f}. \quad (\text{A.36c})$$

Note that, given β_i 's in (23), we have $\mathbf{D}_{\sigma_f} - \frac{1}{2}\mathbf{D}_\rho \succ \mathbf{0}$ and $\mathbf{G} - \frac{1}{\rho}\tilde{\mathbf{A}}^T \mathbf{D}_{L_f} \mathbf{D}_\rho^{-1} \tilde{\mathbf{A}} \succ \mathbf{0}$ (see (A.20) and (A.21)); moreover, since \mathbf{A}_i 's are full column rank, we have $\mathbf{M} \succ \mathbf{0}$. Hence there must exist some $\delta > 0$ such that the three conditions in (A.36) all hold true. \blacksquare

APPENDIX B

PROOF OF THEOREM 3

Proof of Theorem 3(a): Let $\mathbf{x}^* \triangleq [(\mathbf{x}_1^*)^T, \dots, (\mathbf{x}_N^*)^T]^T$ and $\boldsymbol{\nu}^*$ be a pair of optimal primal and dual solutions to (P2), and let $\tilde{\boldsymbol{\nu}}^* \triangleq [(\boldsymbol{\nu}_1^*)^T, \dots, (\boldsymbol{\nu}_N^*)^T]^T$ and $\{\mathbf{u}_{ij}^*, \mathbf{v}_{ij}^*, j \in \mathcal{N}_i\}_{i=1}^N$ be a pair of optimal primal and dual solutions to problem (26). Then they respectively satisfy the following optimality conditions

$$\mathbf{A}_i^T \nabla f_i(\mathbf{x}_i^*) + \partial g_i(\mathbf{x}_i^*) + \mathbf{E}_i^T \boldsymbol{\nu}^* = \mathbf{0}, i \in V, \quad (\text{A.37})$$

$$\sum_{i=1}^N \mathbf{E}_i \mathbf{x}_i^* = \mathbf{q}, \quad (\text{A.38})$$

$$\partial \varphi_i(\boldsymbol{\nu}_i^*) + \frac{1}{N}\mathbf{q} + \sum_{j \in \mathcal{N}_i} (\mathbf{u}_{ij}^* + \mathbf{v}_{ji}^*) = \mathbf{0}, i \in V, \quad (\text{A.39})$$

$$\boldsymbol{\nu}_i^* = \boldsymbol{\nu}_j^* \forall j \in \mathcal{N}_i, i \in V, \quad (\text{A.40})$$

$$\mathbf{u}_{ij}^* + \mathbf{v}_{ij}^* = \mathbf{0} \forall j \in \mathcal{N}_i, i \in V. \quad (\text{A.41})$$

where $\partial \varphi_i(\boldsymbol{\nu}_i^*) = -\mathbf{E}_i \mathbf{x}_i^*$ as \mathbf{x}_i^* is a maximizer to (6) with $\boldsymbol{\nu} = \boldsymbol{\nu}_i^*$ [42]. Under Assumption 1, (A.2) implies that $\boldsymbol{\nu}^* \triangleq \boldsymbol{\nu}_1^* = \dots = \boldsymbol{\nu}_N^*$ and $\tilde{\boldsymbol{\nu}}^* = \mathbf{1}_N \otimes \boldsymbol{\nu}^*$.

Firstly, by recalling that $\mathbf{p}_i^{(k)} = \sum_{j \in \mathcal{N}_i} (\mathbf{u}_{ij}^{(k)} + \mathbf{v}_{ji}^{(k)})$, it follows from (41) and (A.39) that

$$\begin{aligned} \mathbf{0} = & -(\mathbf{E}_i \mathbf{x}_i^{(k)} - \frac{1}{N} \mathbf{q}) + \sum_{j \in \mathcal{N}_i} (\mathbf{u}_{ij}^{(k+1)} + \mathbf{v}_{ji}^{(k+1)}) \\ & + c \sum_{j \in \mathcal{N}_i} (\boldsymbol{\nu}_i^{(k)} + \boldsymbol{\nu}_j^{(k)} - \boldsymbol{\nu}_i^{(k-1)} - \boldsymbol{\nu}_j^{(k-1)}) \end{aligned} \quad (\text{A.42})$$

$$= -\mathbf{E}_i \mathbf{x}_i^* + \frac{1}{N} \mathbf{q} + \sum_{j \in \mathcal{N}_i} \mathbf{u}_{ij}^* + \sum_{j \in \mathcal{N}_i} \mathbf{v}_{ji}^*. \quad (\text{A.43})$$

By multiplying $\boldsymbol{\nu}_i^{(k)} - \boldsymbol{\nu}^*$ to the both sides of (A.43), we obtain

$$\begin{aligned} & \sum_{j \in \mathcal{N}_i} (\mathbf{u}_{ij}^{(k+1)} + \mathbf{v}_{ji}^{(k+1)} - \mathbf{u}_{ij}^* - \mathbf{v}_{ji}^*)^T (\boldsymbol{\nu}_i^{(k)} - \boldsymbol{\nu}^*) \\ & + c \sum_{j \in \mathcal{N}_i} (\boldsymbol{\nu}_i^{(k)} + \boldsymbol{\nu}_j^{(k)} - \boldsymbol{\nu}_i^{(k-1)} - \boldsymbol{\nu}_j^{(k-1)})^T (\boldsymbol{\nu}_i^{(k)} - \boldsymbol{\nu}^*) \\ & - (\mathbf{x}_i^{(k)} - \mathbf{x}_i^*)^T \mathbf{E}_i^T (\boldsymbol{\nu}_i^{(k)} - \boldsymbol{\nu}^*) = \mathbf{0}. \end{aligned} \quad (\text{A.44})$$

Secondly, from the optimality of (43), we have that

$$\begin{aligned} \mathbf{0} = & \mathbf{A}_i^T \nabla f_i(\mathbf{A}_i \mathbf{x}_i^{(k-1)}) + \partial g(\mathbf{x}_i^{(k)}) \\ & + \frac{1}{2|\mathcal{N}_i|} \mathbf{E}_i^T \left[\frac{1}{c} (\mathbf{E}_i \mathbf{x}_i^{(k)} - \frac{1}{N} \mathbf{q}) - \frac{1}{c} \sum_{j \in \mathcal{N}_i} (\mathbf{u}_{ij}^{(k)} + \mathbf{v}_{ji}^{(k)}) \right. \\ & \left. + \sum_{j \in \mathcal{N}_i} (\boldsymbol{\nu}_i^{(k-1)} + \boldsymbol{\nu}_j^{(k-1)}) \right] + \mathbf{P}_i (\mathbf{x}_i^{(k)} - \mathbf{x}_i^{(k-1)}) \\ = & \mathbf{A}_i^T \nabla f_i(\mathbf{A}_i \mathbf{x}_i^{(k-1)}) + \partial g(\mathbf{x}_i^{(k)}) + \mathbf{E}_i^T \boldsymbol{\nu}_i^{(k)} \\ & + \mathbf{P}_i (\mathbf{x}_i^{(k)} - \mathbf{x}_i^{(k-1)}) \end{aligned} \quad (\text{A.45})$$

$$\begin{aligned} = & \mathbf{A}_i^T (\nabla f_i(\mathbf{A}_i \mathbf{x}_i^{(k-1)}) - \nabla f_i(\mathbf{A}_i \mathbf{x}_i^{(k)})) + \mathbf{A}_i^T \nabla f_i(\mathbf{A}_i \mathbf{x}_i^{(k)}) \\ & + \partial g(\mathbf{x}_i^{(k)}) + \mathbf{E}_i^T \boldsymbol{\nu}_i^{(k)} + \mathbf{P}_i (\mathbf{x}_i^{(k)} - \mathbf{x}_i^{(k-1)}) \end{aligned} \quad (\text{A.46})$$

$$= \mathbf{A}_i^T \nabla f_i(\mathbf{A}_i \mathbf{x}_i^*) + \partial g(\mathbf{x}_i^*) + \mathbf{E}_i^T \boldsymbol{\nu}^*, \quad (\text{A.47})$$

where, in the first equality, we have added and subtracted $\frac{1}{2c|\mathcal{N}_i|} \mathbf{E}_i^T \mathbf{E}_i \mathbf{x}_i^{(k)}$ and defined

$$\mathbf{P}_i \triangleq \beta_i \mathbf{I}_K - \frac{1}{2c|\mathcal{N}_i|} \mathbf{E}_i^T \mathbf{E}_i; \quad (\text{A.48})$$

the second equality is due to (33); and the last equality is because \mathbf{x}_i^* is a maximizer to (6) with $\boldsymbol{\nu} = \boldsymbol{\nu}_i^*$. Multiplying both (A.46) and (A.47) with $\mathbf{x}_i^{(k)} - \mathbf{x}_i^*$, combining with (A.44), and summing for

$i = 1, \dots, N$, yields

$$\begin{aligned}
& \sum_{i=1}^N (\nabla f_i(\mathbf{A}_i \mathbf{x}_i^{(k-1)}) - \nabla f_i(\mathbf{A}_i \mathbf{x}_i^{(k)}))^T \mathbf{A}_i (\mathbf{x}_i^{(k)} - \mathbf{x}_i^*) + \sum_{i=1}^N (\mathbf{x}_i^{(k)} - \mathbf{x}_i^{(k-1)})^T \mathbf{P}_i (\mathbf{x}_i^{(k)} - \mathbf{x}_i^*) + \\
& \sum_{i=1}^N (\nabla f_i(\mathbf{A}_i \mathbf{x}_i^{(k)}) - \nabla f_i(\mathbf{A}_i \mathbf{x}_i^*))^T \mathbf{A}_i (\mathbf{x}_i^{(k)} - \mathbf{x}_i^*) + \sum_{i=1}^N (\partial g(\mathbf{x}_i^{(k)}) - \partial g(\mathbf{x}_i^*))^T (\mathbf{x}_i^{(k)} - \mathbf{x}_i^*) \\
& + \sum_{i=1}^N \sum_{j \in \mathcal{N}_i} (\mathbf{u}_{ij}^{(k+1)} + \mathbf{v}_{ji}^{(k+1)} - \mathbf{u}_{ij}^* - \mathbf{v}_{ji}^*)^T (\boldsymbol{\nu}_i^{(k)} - \boldsymbol{\nu}_i^*) \\
& + c \sum_{i=1}^N \sum_{j \in \mathcal{N}_i} (\boldsymbol{\nu}_i^{(k)} + \boldsymbol{\nu}_j^{(k)} - \boldsymbol{\nu}_i^{(k-1)} - \boldsymbol{\nu}_j^{(k-1)})^T (\boldsymbol{\nu}_i^{(k)} - \boldsymbol{\nu}_i^*) = \mathbf{0}.
\end{aligned} \tag{A.49}$$

Similar to (A.15) and by (29), the fifth term in the LHS of (A.49) can be expressed as

$$\begin{aligned}
& \sum_{i=1}^N \sum_{j \in \mathcal{N}_i} (\mathbf{u}_{ij}^{(k+1)} + \mathbf{v}_{ji}^{(k+1)} - \mathbf{u}_{ij}^* - \mathbf{v}_{ji}^*)^T (\boldsymbol{\nu}_i^{(k)} - \boldsymbol{\nu}_i^*) \\
& = \frac{2}{c} (\mathbf{u}^{(k+1)} - \mathbf{u}^*)^T (\mathbf{u}^{(k+1)} - \mathbf{u}^{(k)}).
\end{aligned} \tag{A.50}$$

Moreover, the sixth term in the LHS of (A.49) can be shown as

$$\begin{aligned}
& c \sum_{i=1}^N \sum_{j \in \mathcal{N}_i} (\boldsymbol{\nu}_i^{(k)} - \boldsymbol{\nu}_i^{(k-1)})^T (\boldsymbol{\nu}_i^{(k)} - \boldsymbol{\nu}_i^*) \\
& + c \sum_{i=1}^N \sum_{j \in \mathcal{N}_i} (\boldsymbol{\nu}_j^{(k)} - \boldsymbol{\nu}_j^{(k-1)})^T (\boldsymbol{\nu}_i^{(k)} - \boldsymbol{\nu}_i^*) \\
& = c \sum_{i=1}^N |\mathcal{N}_i| (\boldsymbol{\nu}_i^{(k)} - \boldsymbol{\nu}_i^{(k-1)})^T (\boldsymbol{\nu}_i^{(k)} - \boldsymbol{\nu}_i^*) \\
& + c \sum_{i=1}^N \sum_{j=1}^N [\mathbf{W}]_{i,j} (\boldsymbol{\nu}_j^{(k)} - \boldsymbol{\nu}_j^{(k-1)})^T (\boldsymbol{\nu}_i^{(k)} - \boldsymbol{\nu}_i^*) \\
& = c (\boldsymbol{\nu}^{(k)} - \boldsymbol{\nu}^{(k-1)})^T \mathbf{Q} (\boldsymbol{\nu}^{(k)} - \tilde{\boldsymbol{\nu}}^*),
\end{aligned} \tag{A.51}$$

where $\mathbf{Q} \triangleq (\mathbf{D} + \mathbf{W}) \otimes \mathbf{I}_M$. By applying (A.7), (A.8), (A.9), (A.50) and (A.51) to (A.49), one obtains

$$\begin{aligned}
& \|\mathbf{x}^{(k)} - \mathbf{x}^*\|_{\mathbf{M}}^2 - \frac{1}{2} \|\mathbf{x}^{(k-1)} - \mathbf{x}^{(k)}\|_{\tilde{\mathbf{A}}^T \mathbf{D}_{L_f} \mathbf{D}_{\rho}^{-1} \tilde{\mathbf{A}}}^2 \\
& + (\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)})^T \mathbf{P} (\mathbf{x}^{(k)} - \mathbf{x}^*) \\
& + c (\boldsymbol{\nu}^{(k)} - \boldsymbol{\nu}^{(k-1)})^T \mathbf{Q} (\boldsymbol{\nu}^{(k)} - \tilde{\boldsymbol{\nu}}^*) \\
& + \frac{2}{c} (\mathbf{u}^{(k+1)} - \mathbf{u}^*)^T (\mathbf{u}^{(k+1)} - \mathbf{u}^{(k)}) \leq 0,
\end{aligned} \tag{A.52}$$

where $\boldsymbol{\nu}^{(k)} = [(\boldsymbol{\nu}_1^{(k)})^T, \dots, (\boldsymbol{\nu}_N^{(k)})^T]^T$, $\mathbf{P} = \text{blkdiag}\{\mathbf{P}_1, \dots, \mathbf{P}_N\} \succ \mathbf{0}$, and \mathbf{D}_{σ_f} , \mathbf{D}_{L_f} , \mathbf{D}_{ρ} , $\tilde{\mathbf{A}}$ and \mathbf{M}

are all defined below (A.11). After applying (A.17) to (A.52), we obtain

$$\begin{aligned}
& \|\mathbf{x}^{(k)} - \mathbf{x}^*\|_{\mathbf{M} + \frac{1}{2}\mathbf{P}}^2 + \frac{1}{c}\|\mathbf{u}^{(k+1)} - \mathbf{u}^*\|_2^2 + \frac{c}{2}\|\boldsymbol{\nu}^{(k)} - \tilde{\boldsymbol{\nu}}^*\|_Q^2 \\
& \leq \|\mathbf{x}^{(k-1)} - \mathbf{x}^*\|_{\frac{1}{2}\mathbf{P}}^2 + \frac{1}{c}\|\mathbf{u}^{(k)} - \mathbf{u}^*\|_2^2 + \frac{c}{2}\|\boldsymbol{\nu}^{(k-1)} - \tilde{\boldsymbol{\nu}}^*\|_Q^2 \\
& \quad - \frac{1}{2}\|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\|_{\mathbf{P} - \tilde{\mathbf{A}}^T \mathbf{D}_{L_f} \mathbf{D}_\rho^{-1} \tilde{\mathbf{A}}}^2 - \frac{1}{c}\|\mathbf{u}^{(k+1)} - \mathbf{u}^{(k)}\|_2^2 \\
& \quad - \frac{c}{2}\|\boldsymbol{\nu}^{(k)} - \boldsymbol{\nu}^{(k-1)}\|_Q^2.
\end{aligned} \tag{A.53}$$

It is easy to show that, under (47), it holds true that

$$\sigma_{f,i}^2 - \frac{\rho_i}{2} > 0, \quad \mathbf{P}_i - \frac{L_{f,i}^2}{\rho_i} \mathbf{A}_i^T \mathbf{A}_i \succ \mathbf{0}, \quad \forall i \in V, \tag{A.54}$$

for some $\sigma_{f,i}^2 \leq \rho_i < 2\sigma_{f,i}^2 \forall i \in V$, which implies that

$$\mathbf{P} \succ \mathbf{0}, \quad \mathbf{P} - \tilde{\mathbf{A}}^T \mathbf{D}_{L_f} \mathbf{D}_\rho^{-1} \tilde{\mathbf{A}} \succ \mathbf{0}.$$

Thus, (A.53) implies that **(R1)** the sequence $\|\mathbf{x}^{(k)} - \mathbf{x}^*\|_{\mathbf{M} + \frac{1}{2}\mathbf{P}}^2 + \frac{1}{c}\|\mathbf{u}^{(k+1)} - \mathbf{u}^*\|_2^2 + \frac{c}{2}\|\boldsymbol{\nu}^{(k)} - \tilde{\boldsymbol{\nu}}^*\|_Q^2$ converges for any optimal \mathbf{x}^* to **(P2)**, and optimal $\tilde{\boldsymbol{\nu}}^*$ and \mathbf{u}^* to problem (26); and **(R2)**

$$\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)} \rightarrow \mathbf{0}, \quad \mathbf{u}^{(k+1)} - \mathbf{u}^{(k)} \rightarrow \mathbf{0}, \tag{A.55}$$

$$\|\boldsymbol{\nu}^{(k)} - \boldsymbol{\nu}^{(k-1)}\|_Q^2 \rightarrow 0. \tag{A.56}$$

Let $\hat{\mathbf{x}} = [(\hat{\mathbf{x}}_1)^T, \dots, (\hat{\mathbf{x}}_N)^T]^T$, $\tilde{\boldsymbol{\nu}} = [(\hat{\boldsymbol{\nu}}_1)^T, \dots, (\hat{\boldsymbol{\nu}}_N)^T]^T$, $\hat{\mathbf{u}}_{ij}$ and $\hat{\mathbf{v}}_{ij}$ be a set of limit points of $\{\mathbf{x}^{(k)}\}$, $\{\boldsymbol{\nu}_1^{(k)}, \dots, \boldsymbol{\nu}_N^{(k)}\}$, $\{\mathbf{u}_{ij}^{(k)}\}$ and $\{\mathbf{v}_{ij}^{(k)}\}$, respectively. Firstly, by applying the fact of $\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)} \rightarrow \mathbf{0}$ to (A.46), we have

$$\mathbf{A}_i^T \nabla f_i(\mathbf{A}_i \hat{\mathbf{x}}_i) + \partial g(\hat{\mathbf{x}}_i) + \mathbf{E}_i^T \hat{\boldsymbol{\nu}}_i = \mathbf{0}, \quad \forall i \in V. \tag{A.57}$$

Secondly, by (A.13), we have

$$\hat{\mathbf{u}}_{ij} + \hat{\mathbf{v}}_{ij} = \mathbf{0} \quad \forall j, i. \tag{A.58}$$

Thirdly, applying the fact of $\mathbf{u}_{ij}^{(k+1)} - \mathbf{u}_{ij}^{(k)} \rightarrow \mathbf{0}$ to (29a) yields

$$\boldsymbol{\nu}_i^{(k)} - \boldsymbol{\nu}_j^{(k)} \rightarrow \mathbf{0} \implies \hat{\boldsymbol{\nu}} \triangleq \hat{\boldsymbol{\nu}}_i = \hat{\boldsymbol{\nu}}_j \quad \forall j \in \mathcal{N}_i, i \in V \tag{A.59}$$

The result of $\boldsymbol{\nu}_i^{(k)} - \boldsymbol{\nu}_j^{(k)} \rightarrow \mathbf{0} \forall j, i$ and Assumption 1 implies that $\boldsymbol{\nu}^{(k)} - \mathbf{1}_N \otimes \boldsymbol{\nu}_i^{(k)} \rightarrow \mathbf{0}$ for any $i \in V$. Since the Laplacian matrix $\mathbf{L}\mathbf{1}_N = \mathbf{0}$ [32], one obtains

$$\begin{aligned}
& \|\boldsymbol{\nu}^{(k)} - \boldsymbol{\nu}^{(k-1)}\|_Q^2 \\
& \rightarrow (\mathbf{1}_N \otimes (\boldsymbol{\nu}_i^{(k)} - \boldsymbol{\nu}_i^{(k-1)}))^T \mathbf{Q} (\mathbf{1}_N \otimes (\boldsymbol{\nu}_i^{(k)} - \boldsymbol{\nu}_i^{(k-1)})) \\
& = (\mathbf{1}_N^T (\mathbf{D} + \mathbf{W}) \mathbf{1}_N) \|\boldsymbol{\nu}_i^{(k)} - \boldsymbol{\nu}_i^{(k-1)}\|_2^2 \\
& = (\mathbf{1}_N^T (2\mathbf{D} - \mathbf{L}) \mathbf{1}_N) \|\boldsymbol{\nu}_i^{(k)} - \boldsymbol{\nu}_i^{(k-1)}\|_2^2 \\
& = (2 \sum_{j=1}^N |\mathcal{N}_j|) \|\boldsymbol{\nu}_i^{(k)} - \boldsymbol{\nu}_i^{(k-1)}\|_2^2,
\end{aligned} \tag{A.60}$$

which, when combined with (A.56), further implies that

$$\boldsymbol{\nu}_i^{(k)} - \boldsymbol{\nu}_i^{(k-1)} \rightarrow \mathbf{0} \forall i \in V. \tag{A.61}$$

By applying (A.61) to (A.42), one obtains

$$\mathbf{0} = -\mathbf{E}_i \hat{\mathbf{x}}_i + \frac{1}{N} \mathbf{q} + \sum_{j \in \mathcal{N}_i} \hat{\mathbf{u}}_{ij} + \sum_{j \in \mathcal{N}_i} \hat{\mathbf{v}}_{ji} \tag{A.62}$$

$$= \partial \varphi_i(\hat{\boldsymbol{\nu}}_i) + \frac{1}{N} \mathbf{q} + \sum_{j \in \mathcal{N}_i} \hat{\mathbf{u}}_{ij} + \sum_{j \in \mathcal{N}_i} \hat{\mathbf{v}}_{ji}, \tag{A.63}$$

where $\partial \varphi_i(\hat{\boldsymbol{\nu}}_i) = -\mathbf{E}_i \hat{\mathbf{x}}_i$ since (A.57) implies that $\hat{\mathbf{x}}_i$ is a maximizer to (6) with $\boldsymbol{\nu} = \hat{\boldsymbol{\nu}}_i$ [42]. Finally, by summing (A.62) for $i = 1, \dots, N$, followed by applying (A.14) and (A.58), one obtains

$$\sum_{i=1}^N \mathbf{E}_i \hat{\mathbf{x}}_i = \mathbf{q}. \tag{A.64}$$

The results in (A.57), (A.58), (A.59), (A.63) and (A.64) imply that $\hat{\mathbf{x}}$ and $\hat{\boldsymbol{\nu}}$ are in fact a pair of optimal primal and dual solutions to (P2), and $\tilde{\boldsymbol{\nu}}$ and $\{\hat{\mathbf{u}}_{ij}, \hat{\mathbf{v}}_{ij}\}$ are a pair of optimal primal and dual solutions to problem (26) [see (A.37) to (A.41)]. Thus, according to (R1), the sequence $\|\mathbf{x}^{(k)} - \hat{\mathbf{x}}\|_{\mathbf{M} + \frac{1}{2}\mathbf{P}}^2 + \frac{1}{c} \|\mathbf{u}^{(k+1)} - \hat{\mathbf{u}}\|_2^2 + \frac{c}{2} \|\boldsymbol{\nu}^{(k)} - \tilde{\boldsymbol{\nu}}\|_Q^2$ in fact converges to zero and thereby $\mathbf{x}^{(k)} \rightarrow \hat{\mathbf{x}}$, $\mathbf{u}^{(k+1)} \rightarrow \hat{\mathbf{u}}$ and $\boldsymbol{\nu}_i^{(k)} \rightarrow \hat{\boldsymbol{\nu}} \forall i \in V$. ■

Proof of Theorem 3(b): We assume that $\phi_i(\mathbf{x}_i) = f_i(\mathbf{A}_i \mathbf{x}_i)$, \mathbf{A}_i has full column rank and \mathbf{E}_i has full row rank, for all $i \in V$. Denote $\mathbf{r}^{(k)} \triangleq \|\mathbf{x}^{(k)} - \mathbf{x}^*\|_{\alpha\mathbf{M} + \frac{1}{2}\mathbf{P}}^2 + \frac{1}{c} \|\mathbf{u}^{(k+1)} - \mathbf{u}^*\|_2^2 + \frac{c}{2} \|\boldsymbol{\nu}^{(k)} - \tilde{\boldsymbol{\nu}}^*\|_Q^2$ for some $\alpha > 0$. One can write (A.53) as follows

$$\begin{aligned}
& \mathbf{r}^{(k)} + \|\mathbf{x}^{(k)} - \mathbf{x}^*\|_{(1-\alpha)\mathbf{M}}^2 + \|\mathbf{x}^{(k-1)} - \mathbf{x}^*\|_{\alpha\mathbf{M}}^2 \\
& \leq \mathbf{r}^{(k-1)} - \frac{1}{2} \|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\|_{\mathbf{P} - \frac{1}{2}\bar{\mathbf{A}}^T \mathbf{D}_{L_f} \mathbf{D}_\rho^{-1} \bar{\mathbf{A}}}^2 \\
& \quad - \frac{1}{c} \|\mathbf{u}^{(k+1)} - \mathbf{u}^{(k)}\|_2^2 - \frac{c}{2} \|\boldsymbol{\nu}^{(k)} - \boldsymbol{\nu}^{(k-1)}\|_Q^2.
\end{aligned}$$

Therefore, it suffices to show that, for some $\delta > 0$,

$$\begin{aligned} & \|\mathbf{x}^{(k)} - \mathbf{x}^*\|_{(1-\alpha)\mathbf{M}}^2 + \|\mathbf{x}^{(k-1)} - \mathbf{x}^*\|_{\alpha\mathbf{M}}^2 \\ & + \frac{1}{2}\|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\|_{\mathbf{P} - \frac{1}{2}\tilde{\mathbf{A}}^T \mathbf{D}_{L_f} \mathbf{D}_\rho^{-1} \tilde{\mathbf{A}}}^2 + \frac{1}{c}\|\mathbf{u}^{(k+1)} - \mathbf{u}^{(k)}\|_2^2 \\ & + \frac{c}{2}\|\boldsymbol{\nu}^{(k)} - \boldsymbol{\nu}^{(k-1)}\|_{\mathbf{Q}}^2 \geq \delta \mathbf{r}^{(k)}. \end{aligned} \quad (\text{A.65})$$

Firstly, from (A.45) and (A.47), we have that (without g_i 's)

$$\begin{aligned} & \mathbf{A}_i^T (\nabla f_i(\mathbf{A}_i \mathbf{x}_i^{(k-1)}) - \nabla f_i(\mathbf{A}_i \mathbf{x}_i^*)) + \mathbf{E}_i^T (\boldsymbol{\nu}_i^{(k)} - \boldsymbol{\nu}^*) \\ & + \mathbf{P}_i (\mathbf{x}_i^{(k)} - \mathbf{x}_i^{(k-1)}) = \mathbf{0}. \end{aligned} \quad (\text{A.66})$$

By applying (A.34) to (A.66), we have, for some $\mu_1 > 1$,

$$\begin{aligned} & \|\mathbf{P}_i (\mathbf{x}_i^{(k)} - \mathbf{x}_i^{(k-1)})\|^2 \\ & \geq (1 - \mu_1) \|\mathbf{A}_i^T (\nabla f_i(\mathbf{A}_i \mathbf{x}_i^{(k-1)}) - \nabla f_i(\mathbf{A}_i \mathbf{x}_i^*))\|^2 \\ & \quad + (1 - \frac{1}{\mu_1}) \|\mathbf{E}_i^T (\boldsymbol{\nu}_i^{(k)} - \boldsymbol{\nu}^*)\|_2^2 \\ & \geq (1 - \mu_1) L_{f,i} \lambda_{\max}^2(\mathbf{A}_i^T \mathbf{A}_i) \|\mathbf{x}_i^{(k-1)} - \mathbf{x}_i^*\|_2^2 \\ & \quad + (1 - \frac{1}{\mu_1}) \lambda_{\min}(\mathbf{E}_i \mathbf{E}_i^T) \|\boldsymbol{\nu}_i^{(k)} - \boldsymbol{\nu}^*\|_2^2, \end{aligned} \quad (\text{A.67})$$

where the second inequality is obtained by (12). Note that $\mathbf{D} + \mathbf{W} = 2\mathbf{D} - \mathbf{L} \preceq 2\mathbf{D}$ as $\mathbf{L} \succeq \mathbf{0}$ [32].

Hence, we have

$$\begin{aligned} & \frac{c\delta}{2} \|\boldsymbol{\nu}^{(k)} - \tilde{\boldsymbol{\nu}}^*\|_{\mathbf{Q}}^2 \leq c\delta \|\boldsymbol{\nu}^{(k)} - \tilde{\boldsymbol{\nu}}^*\|_{\mathbf{D} \otimes \mathbf{I}_M}^2 \\ & \leq c\delta \tau_1 \|(\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)})\|_{\mathbf{P}^T \mathbf{P}}^2 + c\delta \tau_2 \|\mathbf{x}^{(k-1)} - \mathbf{x}^*\|_2^2, \end{aligned} \quad (\text{A.68})$$

where the second inequality is due to (A.67), $\tau_1 = \max_{i \in V} \left\{ \frac{|\mathcal{N}_i|}{(1 - \frac{1}{\mu_1}) \lambda_{\min}(\mathbf{E}_i \mathbf{E}_i^T)} \right\} > 0$ and $\tau_2 = \max_{i \in V} \left\{ \frac{(\mu_1 - 1) \lambda_{\max}^2(\mathbf{A}_i^T \mathbf{A}_i) |\mathcal{N}_i|}{(1 - \frac{1}{\mu_1}) \lambda_{\min}(\mathbf{E}_i \mathbf{E}_i^T)} \right\}$ are finite given that \mathbf{E}_i 's have full row rank.

Secondly, upon stacking (A.43) for all $i \in V$ and applying (A.3) and (A.12), one obtains

$$\begin{aligned} & \boldsymbol{\Upsilon} (\mathbf{u}^{(k+1)} - \mathbf{u}^*) + c\mathbf{Q} (\boldsymbol{\nu}^{(k)} - \boldsymbol{\nu}^{(k-1)}) \\ & - \mathbf{E} (\mathbf{x}^{(k)} - \mathbf{x}^*) = \mathbf{0}, \end{aligned} \quad (\text{A.69})$$

where $\mathbf{E} = \text{blkdiag}\{\mathbf{E}_1, \dots, \mathbf{E}_N\}$ and $\boldsymbol{\Upsilon}$ is given in (A.31). Analogously, by applying (A.34) to (A.69)

and by (A.32), one can show that, for some $\mu_2 > 1$,

$$\begin{aligned} \frac{\delta}{c} \|\mathbf{u}^{(k+1)} - \mathbf{u}^*\|_2^2 &\leq \frac{\delta}{c\tau_3} \|\mathbf{x}^{(k)} - \mathbf{x}^*\|_{\mathbf{E}^T \mathbf{E}}^2 \\ &\quad + \frac{\delta(\mu_2 - 1)c}{\tau_3} \|\boldsymbol{\nu}^{(k)} - \boldsymbol{\nu}^{(k-1)}\|_{\mathbf{Q}}^2, \end{aligned} \quad (\text{A.70})$$

where $\tau_3 = (1 - \frac{1}{\mu_2})\sigma_{\min}^2(\boldsymbol{\Upsilon}) > 0$. By (A.68) and (A.70), sufficient conditions for satisfying (A.65) are therefore given by: $\forall i \in V$,

$$(1 - \alpha - \delta\alpha)(\sigma_{f,i}^2 - \frac{\rho_i}{2})\mathbf{A}_i^T \mathbf{A}_i \succeq \frac{\delta}{2}\mathbf{P}_i + \frac{\delta}{c\tau_3}\mathbf{A}_i^T \mathbf{A}_i, \quad (\text{A.71a})$$

$$\alpha(\sigma_{f,i}^2 - \frac{\rho_i}{2})\mathbf{A}_i^T \mathbf{A}_i \succeq c\delta\tau_2 \mathbf{I}_K, \quad (\text{A.71b})$$

$$\frac{1}{2}\mathbf{P}_i - \frac{L_{f,i}^2}{2\rho_i}\mathbf{A}_i^T \mathbf{A}_i \succeq c\delta\tau_1 \mathbf{P}_i^T \mathbf{P}_i, \quad (\text{A.71c})$$

$$\frac{1}{2} \geq \frac{\delta(\mu_2 - 1)}{\tau_3}. \quad (\text{A.71d})$$

Under (A.54) and full column rank \mathbf{A}_i 's, we see that (A.71) is true for some $\delta > 0$. The proof is complete. \blacksquare

REFERENCES

- [1] I. Foster, Y. Zhao, I. Raicu, and S. Lu, "Cloud computing and grid computing 360-degree compared," in *Proc. Grid Computing Environments Workshop*, Austin, TX, USA, Nov. 12-16, 2008, pp. 1–10.
- [2] R. Bekkerman, M. Bilenko, and J. Langford, *Scaling up Machine Learning- Parallel and Distributed Approaches*. Cambridge University Press, 2012.
- [3] G. R. Andrews, *Foundations of Multithreaded, Parallel, and Distributed Programming*. Addison-Wesley, 2007.
- [4] S. Ghosh, *Distributed Systems- An Algorithmic Approach*. Chapman & Hall/CRC Computer & Information Science Series, 2007.
- [5] A. Nedić, A. Ozdaglar, , and A. Parrilo, "Constrained consensus and optimization in multi-agent networks," *IEEE Trans. Automatic Control*, vol. 55, no. 4, pp. 922–938, April 2010.
- [6] B. Johansson, T. Keviczky, M. Johansson, and K. Johansson, "Subgradient methods and consensus algorithms for solving convex optimization problems," in *Proc. IEEE CDC*, Cancun, Mexico, Dec. 9-11, 2008, pp. 4185–4190.
- [7] M. Zhu and S. Martínez, "On distributed convex optimization under inequality and equality constraints," *IEEE Trans. Automatic Control*, vol. 57, no. 1, pp. 151–164, Jan. 2012.
- [8] J. Chen and A. H. Sayed, "Diffusion adaption strategies for distributed optimization and learning networks," *IEEE. Trans. Signal Process.*, vol. 60, no. 8, pp. 4289–4305, Aug. 2012.
- [9] M. Elad, *Sparse and Redundant Representations*. New York, NY, USA: Springer Science + Business Media, 2010.
- [10] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein, "Distributed optimization and statistical learning via the alternating direction method of multipliers," *Foundations and Trends in Machine Learning*, vol. 3, no. 1, pp. 1–122, 2011.
- [11] S. Chen, D. Donoho, and M. Saunders, "Atomic decomposition by basis pursuit," *SIAM J. Sci. Comput.*, vol. 20, no. 1, pp. 33–61, 1998.
- [12] J. F. C. Mota, J. M. F. Xavier, P. M. Q. Aguiar, and M. Puschel, "Distributed basis pursuit," *IEEE. Trans. Signal Process.*, vol. 60, no. 4, pp. 1942–1956, April 2012.

- [13] D. P. Bertsekas, *Network Optimization : Contribuous and Discrete Models*. Athena Scientific, 1998.
- [14] C. Shen, T.-H. Chang, K.-Y. Wang, Z. Qiu, and C.-Y. Chi, "Distributed robust multicell coordinated beamforming with imperfect CSI: An ADMM approach," *IEEE Trans. Signal Processing*, vol. 60, no. 6, pp. 2988–3003, 2012.
- [15] S. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge, UK: Cambridge University Press, 2004.
- [16] S. S. Ram, A. Nedić, and V. V. Veeravalli, "A new class of distributed optimization algorithm: Application of regression of distributed data," *Optimization Methods and Software*, vol. 27, no. 1, pp. 71–88, 2012.
- [17] T.-H. Chang, A. Nedić, and A. Scaglione, "Distributed constrained optimization by consensus-based primal-dual perturbation method," *IEEE. Trans. Automatic Control.*, vol. 59, no. 6, pp. 1524–1538, June 2014.
- [18] D. P. Bertsekas and J. N. Tsitsiklis, *Parallel and distributed computation: Numerical methods*. Upper Saddle River, NJ, USA: Prentice-Hall, Inc., 1989.
- [19] G. Mateos, J. A. Bazerque, and G. B. Giannakis, "Distributed sparse linear regression," *IEEE Trans. Signal Process.*, vol. 58, no. 10, pp. 5262–5276, Dec. 2010.
- [20] R. Tibshirani, "Regression shrinkage and selection via the LASSO," *J. Roy. Stat. Soc. B*, vol. 58, pp. 267–288, 1996.
- [21] J. F. C. Mota, J. M. F. Xavier, P. M. Q. Aguiar, and M. Puschel, "D-ADMM: A communication-efficient distributed algorithm for separable optimization," *IEEE. Trans. Signal Process.*, vol. 60, no. 10, pp. 2718–2723, May 2013.
- [22] E. Wei and A. Ozdaglar, "On the $O(1/K)$ convergence of asynchronous distributed alternating direction method of multipliers," in *Proc. IEEE GlobalSIP*, Austin, TX, USA, Dec. 3-5, 2013, pp. 551–554.
- [23] —, "Distributed alternating direction method of multipliers," in *Proc. IEEE CDC*, Maui, HI, USA, Dec. 10-13, 2012, pp. 5445–5450.
- [24] W. Deng and W. Yin, "On the global and linear convergence of the generalized alternating direction method of multipliers," Rice CAAM technical report 12-14, 2012.
- [25] M. Hong and Z.-Q. Luo, "On the linear convergence of the alternating direction method of multipliers," available on arxiv.org.
- [26] W. Shi, Q. Ling, K. Yuan, G. Wu, and W. Yin, "On the linear convergence of the ADMM in decentralized consensus optimization," *IEEE Trans. Signal Process.*, vol. 62, no. 7, pp. 1750–1761, April 2014.
- [27] D. Jakovetic, J. M. F. Moura, and J. Xavier, "Linear convergence rate of a class of distributed augmented lagrangian algorithms," available on arxiv.org.
- [28] Y. Nesterov, "Smooth minimization of nonsmooth functions," *Math. Program.*, vol. 103, no. 1, pp. 127–152, 2005.
- [29] B. He and X. Yuan, "Linearized alternating direction method of multipliers with gaussian back substitution for separable convex programming," *Numerical Algebra. Control and Optimization*, vol. 3, no. 2, pp. 247–260, 2013.
- [30] S. Ma, "Alternating proximal gradient method for convex minimization," available on <http://www.optimization-online.org/>.
- [31] I. Foster, Y. Zhao, I. Raicu, and S. Lu, "Large-scale sparse logistic regression," in *Proc. ACM Int. Conf. on Knowledge Discovery and Data Mining*, New York, NY, USA, June 28 - July 1, 2009, pp. 547–556.
- [32] F. R.-K. Chung, *Spectral graph theory*. CBMS Regional Conference Series in Mathematics, No. 92, 1996.
- [33] E. Ghadimi, A. Teixeira, I. Shames, and M. Johansson, "Optimal parameter selection for the alternating direction method of multipliers (ADMM): Quadratic problems," available on arxiv.org.
- [34] P. L. Combettes and J.-C. Pesquet, "Proximal splitting methods in signal processing," available on arxiv.org.
- [35] B. He, Z. Peng, and X. Wang, "Proximal alternating direction-based contraction methods for separable linearly constrained convex optimization," *Frontiers of Math. in China*, vol. 6, no. 1, pp. 79–114, 2011.
- [36] S. Ma and S. Zhang, "An extragradient-based alternating direction method for convex minimization," available on arxiv.org.

- [37] Q. Ling and A. Ribeiro, “Decentralized linearized alternating direction method of multipliers,” in *Proc. IEEE ICASSP*, Florence, Italy, May 4-9, 2014, pp. 5445–5450.
- [38] D. P. Bertsekas, A. Nedić, and A. E. Ozdaglar, *Convex analysis and optimization*. Cambridge, Massachusetts: Athena Scientific, 2003.
- [39] M. E. Yildiz and A. Scaglione, “Coding with side information for rate-constrained consensus,” *IEEE Trans. Signal Process.*, vol. 56, no. 8, pp. 3753–3764, 2008.
- [40] A. Beck and M. Teboulle, “A fast iterative shrinkage-thresholding algorithm for linear inverse problems,” *SIAM J. Imaging Sci.*, vol. 2, no. 1, pp. 183–202, 2009.
- [41] N. Parikh and S. Boyd, “Proximal algorithms,” *Foundations and Trends in Optimization*, vol. 1, no. 3, pp. 1–112, 2013.
- [42] S. Boyd and A. Mutapcic, “Subgradient methods,” available at www.stanford.edu/class/ee392o/subgrad_method.pdf.